

# High School Longitudinal Study of 2009 (HSLs:09) Postsecondary Education Transcript Study and Student Financial Aid Records Collection

Data File Documentation

**AUGUST 2020**

**Michael A. Duprey**  
**Daniel J. Pratt**  
**David H. Wilson**  
**Donna M. Jewell**  
**Derick S. Brown**  
**Lesa R. Caves**  
**Satkartar K. Kinney**  
**Tiffany L. Mattox**  
**Nichole Smith Ritchie**  
**James E. Rogers**  
**Colleen M. Spagnardi**  
**Jamie D. Wescott**  
RTI International

**Elise M. Christopher**  
*Project Officer*  
National Center for Education Statistics

**U.S. Department of Education**

Betsy DeVos

*Secretary***Institute of Education Sciences**

Mark Schneider

*Director***National Center for Education Statistics**

James L. Woodworth

*Commissioner***Sample Surveys Division**

Chris Chapman

*Associate Commissioner*

The National Center for Education Statistics (NCES) is the primary federal entity for collecting, analyzing, and reporting data related to education in the United States and other nations. It fulfills a congressional mandate to collect, report, analyze, and disseminate statistical data related to education in the United States and in other nations; conduct and publish reports and specialized analyses of the meaning and significance of such statistics; assist state and local education agencies in improving their statistical systems; and review and report on education activities in foreign countries.

NCES activities are designed to address high-priority education data needs; provide consistent, reliable, complete, and accurate indicators of education status and trends; and report timely, useful, and high-quality data to the U.S.

Department of Education, the Congress, the states, other education policymakers, practitioners, data users, and the general public. Unless specifically noted, all information contained herein is in the public domain.

We strive to make our products available in a variety of formats and in language that is appropriate to a variety of audiences. You, as our customer, are the best judge of our success in communicating information effectively. If you have any comments or suggestions about this or any other NCES product or report, we would like to hear from you. Please direct your comments to

NCES, IES, U.S. Department of Education  
Potomac Center Plaza (PCP)  
550 12th Street, SW  
Washington, DC 20202

August 2020

The NCES Home Page address is <https://nces.ed.gov>.

The NCES Publications and Products address is <https://nces.ed.gov/pubsearch>.

This report was prepared for the National Center for Education Statistics under Contract No. ED-IES-14-C-0112 with RTI International. Mention of trade names, commercial products, or organizations does not imply endorsement by the U.S. Government.

**Suggested Citation**

Duprey, M.D., Pratt, D.J., Wilson, D.H., Jewell, D.M., Brown, D.S., Caves, L.R., Kinney, S.K., Mattox, T.L., Smith Ritchie, N., Rogers, J.E., Spagnardi, C.M., and Wescott, J.D. (2020). *High School Longitudinal Study of 2009 (HSLs:09) Postsecondary Education Transcript Study and Student Financial Aid Records Collection Data File Documentation* (NCES 2020-004). U.S. Department of Education. Washington, DC: National Center for Education Statistics, Institute of Education Sciences. Retrieved [date] from <https://nces.ed.gov/pubsearch>.

**Content Contact**

National Center for Education Statistics

(800) 677-6987

[NCES.Info@ed.gov](mailto:NCES.Info@ed.gov)

# Contents

	PAGE
Chapter 1. Introduction.....	1
1.1 Overview of Data File Documentation (DFD) Report.....	1
1.2 Historical Background: NCES Secondary Longitudinal Studies Program.....	2
1.3 High School Longitudinal Study of 2009.....	5
1.3.1 Base Year, First Follow-up, 2013 Update, and Second Follow-up.....	6
1.3.2 PETS-SR.....	8
1.3.3 Research and Policy Issues and Analytic Levels .....	9
Chapter 2. Sample Design.....	16
2.1 Base-year Sample Design.....	16
2.2 First Follow-up Sample Design.....	17
2.3 2013 Update and High School Transcript Study Sample Design.....	18
2.4 Second Follow-up Sample Design .....	19
2.5 Postsecondary Education Transcript and Student Records Sample Design .	19
Chapter 3. Data Collection Methodology and Results .....	20
3.1 Postsecondary Education Transcripts and Student Records Systems and Processes.....	20
3.1.1 Postsecondary Data Portal Website.....	20
3.1.2 Institution Contacting Staff Training .....	21
3.1.3 Institution Contacting and Recruitment .....	21
3.2 Postsecondary Education Transcripts-Specific Processes and Quality Control .....	24
3.2.1 Data Receipt Procedures .....	24
3.2.2 Transcript Keying/Coding System and Keyer/Coders .....	25
3.2.3 Coding Taxonomies .....	27
3.2.4 Transcript Keying/Coding Quality Control.....	28
3.3 Postsecondary Education Transcripts Data Collection Results .....	32
3.4 Student Records-Specific Processes and Quality Control.....	34
3.4.1 Student Records Instrument.....	34
3.4.2 Student Records Quality-Control Procedures.....	36
3.5 Student Records Data Collection Results .....	36
Chapter 4. Data Processing and Editing.....	39
4.1 Data Processing .....	39
4.1.1 Transcript Data Reassignment and Consolidation.....	39

	PAGE
4.1.2 Processing Student Records Data.....	40
4.2 Data Editing, Documentation, and Review.....	41
Chapter 5. Response Rates, Analytic Weights, Variance and Design Effects	
Estimation, Nonresponse Bias Analysis, Imputation, and Disclosure	
Avoidance .....	43
5.1 Criteria for Defining Respondents.....	43
5.2 Unit Response Rates .....	44
5.3 Overview of Weighting.....	48
5.3.1 Analysis Weights.....	48
5.3.2 BRR Weights.....	55
5.3.3 Weight Characteristics .....	56
5.3.4 Weighting Quality Control.....	58
5.4 Choosing an Analytic Weight .....	59
5.5 Measures of Precision: Standard Errors and Design Effects.....	70
5.5.1 Standard Errors.....	70
5.5.2 Design Effects .....	74
5.6 Unit and Item Nonresponse Bias Analysis.....	78
5.6.1 Unit Nonresponse Bias Analysis.....	78
5.6.2 Item Nonresponse Bias Analysis .....	82
5.7 Single-value Item Imputation .....	91
5.7.1 Imputed Items .....	92
5.7.2 Evaluation of the Imputed Values.....	95
5.8 Disclosure Risk Analysis and Protections.....	95
5.8.1 PETS-SR Data Products .....	96
5.8.2 Recoding, Suppression, and Swapping.....	96
Chapter 6. Data File Contents.....	98
6.1 PETS-SR Data Products .....	98
6.1.1 Restricted-use Data Products .....	98
6.1.2 Public-use Data Products.....	99
6.2 Contents of the PETS-SR Data Products.....	100
6.3 Variable Naming Schema .....	104
6.4 Missing Data.....	104
6.4.1 Reserve Codes.....	105
6.4.2 Placeholder Records.....	105
6.5 Composite Variables .....	106
6.6 Data Anomalies and Considerations .....	107

## Appendixes

- A. Glossary of Terms
- B. Student Financial Aid Records Instrument Specifications
- C. Notification Materials for Data Collection
- D. Unit and Item Nonresponse Bias Analysis
- E. Detailed Weighting Equations with Specifications
- F. Standard Errors and Design Effects
- G. Imputation Details
- H. ECB Variable Listing
- I. Documentation for Composite Variables

# List of Tables

TABLE	PAGE
1. Upcoding of “other, specify” responses: 2018 .....	31
2. Eligible institution participation, by institution type: 2018 .....	33
3. Student-level transcript collection results: 2018 .....	34
4. Number and percent of participating institutions, by student records collection methods, by institution type: 2018 .....	37
5. Student-level student records collection results: 2018 .....	38
6. HSLS:09 unit response rates .....	47
7. Descriptive characteristics of PETS-SR survey weights: 2018 .....	56
8. Weighted counts and percentages of X2SEX for restricted- and public- use files, by PETS-SR survey weight: 2018 .....	57
9. HSLS:09 analysis weights: 2018 .....	62
10. Number and percentage of completed surveys, high school transcript responses, postsecondary transcript and student records responses, or their combinations for the student sample, and associated recommended weights: PETS-SR .....	65
11. Average design effects ( <i>deff</i> ) and root design effects ( <i>deft</i> ) for postsecondary transcript and student records variables .....	77
12. Summary statistics for unit nonresponse bias analyses before and after weight adjustments for nonresponse, by HSLS:09 PETS-SR analysis weights: 2018 .....	80
13. Student records items with a weighted item response rate below 85 percent using SR student weight (W5PSRECORDS) .....	85
14. Student records items with a weighted item response rate below 85 percent using PETS student weight (W5PSTRANS) .....	87
15. Frequency distribution of the estimated bias ratios for student records items .....	88
16. Frequency distribution of the estimated bias ratios for transcript items .....	89

TABLE	PAGE
17. Summary statistics for student records item nonresponse bias analyses using WSPSRECORDS weight .....	90
18. Summary statistics for student-level item nonresponse bias analyses using W5PSTRANS weight .....	91
19. Student records variables included in single-value imputation, by number and weighted percentage of values missing: 2018 .....	93
20. PETS-SR data products: 2018.....	103
21. Reserve code values: 2018 .....	105

# List of Figures

FIGURE	PAGE
1. Longitudinal design for the NCES secondary longitudinal studies program: 1972–2025.....	4
2. Longitudinal design for the HSLS:09 9th-grade cohort: 2009–2025 .....	6
3. Keying and Coding System Degrees page: 2018 .....	26
4. Keying and Coding System Courses page: 2018 .....	26
5. Code diagram: 2018.....	28
6. PETS-SR unknown eligibility adjustment construction: 2018 .....	52
7. PETS-SR nonresponse and calibration weighting adjustment construction: 2018 .....	54
8. Example SAS-callable SUDAAN code to calculate an estimated mean and linearization standard error for a postsecondary transcript student-level analysis.....	72
9. Example SUDAAN code to calculate an estimated mean and replicate (BRR) standard error for a postsecondary transcript student-level analysis .....	72
10. Example Stata code to calculate an estimated mean and linearization standard error for a postsecondary transcript student-level analysis.....	73
11. Example Stata code to calculate an estimated mean and replicate (BRR) standard error for a postsecondary transcript student-level analysis.....	73
12. Example SAS code to calculate an estimated mean and linearization standard error for a postsecondary transcript student-level analysis.....	73
13. Example SAS code to calculate an estimated mean and replicate (BRR) standard error for a postsecondary transcript student-level analysis.....	73
14. Example R survey package code to calculate an estimated mean and linearization standard error for a postsecondary transcript student-level analysis.....	74
15. Example R survey package code to calculate an estimated mean and replicate (BRR) standard error for a postsecondary transcript student-level analysis.....	74
16. Example IBM SPSS complex samples code to calculate an estimated mean and linearization standard error for a postsecondary transcript student-level analysis .....	74

# Chapter 1. Introduction

## 1.1 Overview of Data File Documentation (DFD) Report

This DFD report provides information and guidance for users of data from the base year through the Postsecondary Education Transcript Study and Student Financial Aid Records Collection (PETS-SR) of the High School Longitudinal Study of 2009 (HSLS:09), with a focus on the PETS-SR data collection. HSLS:09 is sponsored by the National Center for Education Statistics (NCES) of the Institute of Education Sciences, U.S. Department of Education, with additional support from the National Science Foundation.

This documentation is divided into six chapters. Chapter 1 provides an introduction and outlines the organization of the documentation. It describes the historical background of HSLS:09 as part of the NCES secondary longitudinal studies program and supplies a study overview including levels of analysis and research questions. Chapter 1 also briefly describes previous HSLS:09 data collections including surveys with students, parents, and various school personnel; the High School Transcript collection; the survey and administrative data collections that comprise the second follow-up; and the PETS-SR collection components.

Chapter 2 describes the steps used to select the base-year sample and describes sampling through each subsequent follow-up, explaining the resulting sample for the PETS-SR collection. Chapter 2 also describes which sample members are included on the PETS-SR data files.

Chapter 3 discusses the Postsecondary Education Transcript Study (PETS) and Student Financial Aid Records (SR) collection, detailing the collection methodology and results, including data collection design, procedures, participation outcomes, and evaluations.

Chapter 4 discusses the data-processing and post-collection activities for the PETS-SR collection.

Chapter 5 describes response rates, weighting, and other statistical procedures. This chapter presents information on response rates in the base year through the PETS-SR collection. The chapter includes a section explaining the creation of the PETS-SR weights and an overview of the impact of the weights on nonresponse

bias. Also included are sections describing item and unit nonresponse bias analyses, imputation methodology, imputation results, and the disclosure-avoidance procedures applied to the PETS-SR collection data.

Chapter 6 describes the contents of the restricted- and public-use data files from the base year through the PETS-SR collection. The chapter describes how data users can access the restricted- and public-use data, whether through electronic codebook (ECB), or Online Codebook. Variable naming conventions and the scheme used for denoting missing data are also covered. The chapter concludes with a description of the composite variables created from the multiple data sources and analytic weight variables provided in the data files.

This documentation also contains the following appendixes:

- A. Glossary of Terms
- B. Student Financial Aid Records Instrument Specifications
- C. Notification Materials for Data Collection
- D. Unit and Item Nonresponse Bias Analysis
- E. Detailed Weighting Equations with Specifications
- F. Standard Errors and Design Effects
- G. Imputation Details
- H. ECB Variable Listing
- I. Documentation for Composite Variables

## 1.2 Historical Background: NCES Secondary Longitudinal Studies Program

In response to its mandate to “collect, report, analyze, and disseminate statistical data related to education in the United States and in other nations”<sup>1</sup> and the need for policy-relevant, nationally representative longitudinal data on high school students, NCES has maintained a secondary longitudinal studies program. The aim of this continuing program is to study the educational, vocational, and personal development of students at various stages in their educational careers and to examine the personal, familial, social, institutional, and cultural factors that may affect that development.

The program consists of four completed studies, the ongoing HSLS:09, and one new study. The completed studies are the National Longitudinal Study of the High

---

<sup>1</sup> 20 USC § 9543 (a). See <https://www.govinfo.gov/content/pkg/USCODE-2018-title20/pdf/USCODE-2018-title20-chap76.pdf>.

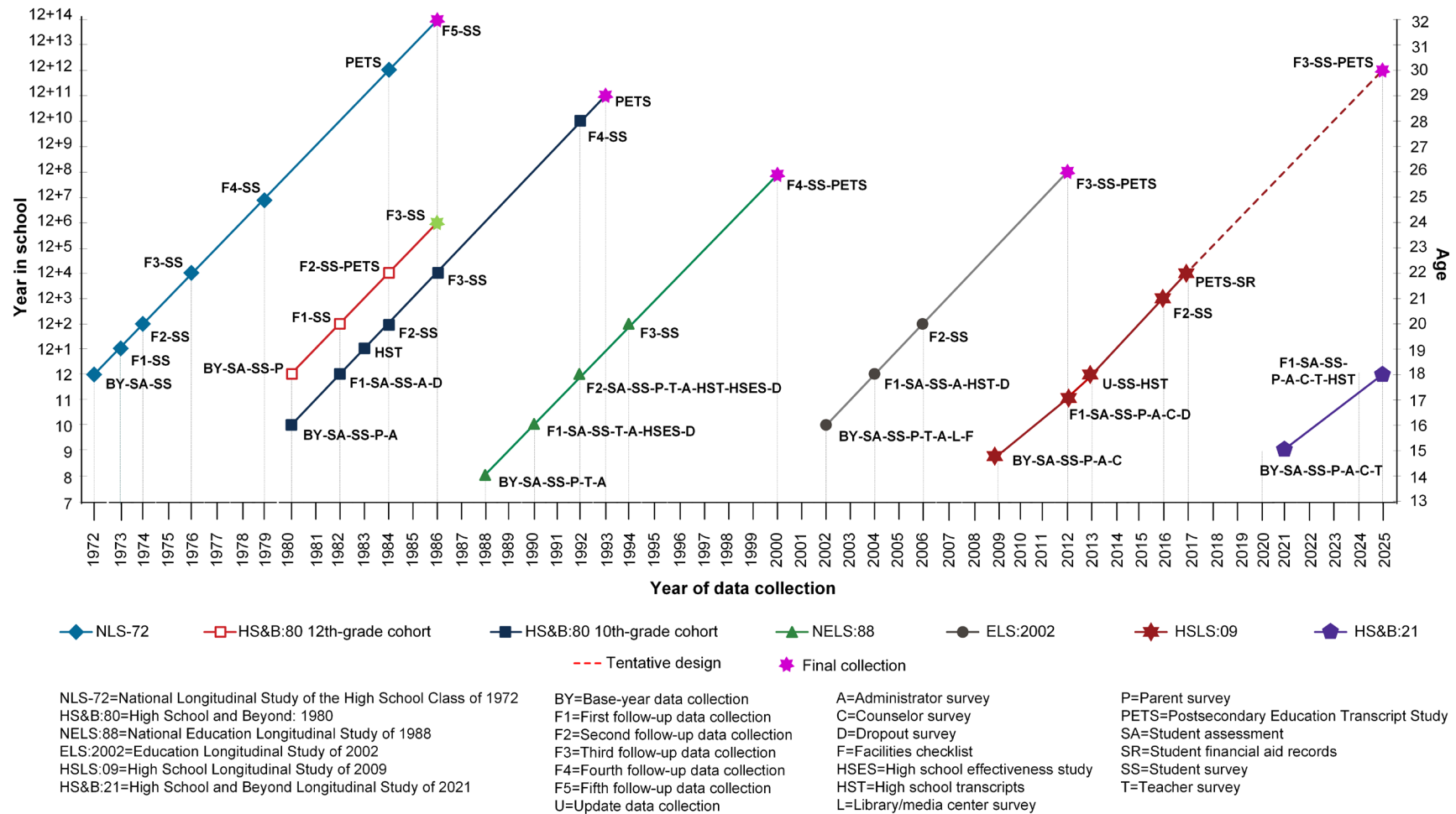
School Class of 1972 (NLS-72), the High School and Beyond Longitudinal Study of 1980 (HS&B:80), the National Education Longitudinal Study of 1988 (NELS:88), and the Education Longitudinal Study of 2002 (ELS:2002). The new study is the High School and Beyond Longitudinal Study of 2021 (HS&B:21), which will begin base-year data collection in the fall of 2021.

Together, these six studies will describe the secondary and postsecondary experiences of students from six decades—the 1970s, 1980s, 1990s, 2000s, 2010s, and 2020s—and provide bases for further understanding the correlates of educational success in the United States. Information on both the current and completed studies in the series is available on the NCES website.<sup>2</sup>

Figure 1 presents a chronology of these six longitudinal education studies and highlights their component and comparison points for the time frame from 1972 to 2025.

---

<sup>2</sup> <https://nces.ed.gov/surveys/slsp/>

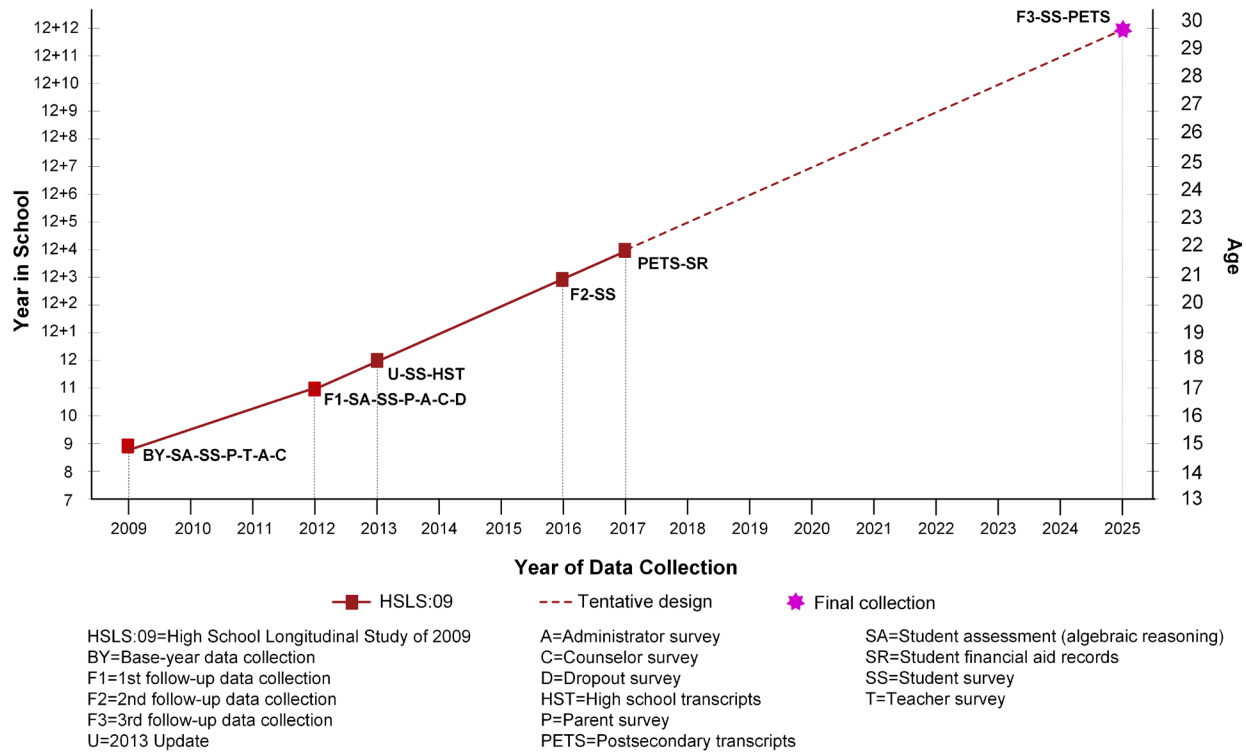
**Figure 1. Longitudinal design for the NCES secondary longitudinal studies program: 1972–2025**

SOURCE: U.S. Department of Education, National Center for Education Statistics, High School Longitudinal Study of 2009 (HSLS:09).

## 1.3 High School Longitudinal Study of 2009

This section provides an overview of the historical context of HSLS:09 in which the PETS-SR collection is situated. The section describes the target population and briefly describes the preceding collections and their respective components.

HSLS:09 is based upon a nationally representative sample of entering 9th-graders in the fall of 2009 who were selected from a nationally representative sample of high schools with 9th and 11th grades. The study is designed to serve multiple policy objectives, primarily through longitudinal analysis. The goal of HSLS:09 is to provide data to understand better the impact of earlier educational experiences, starting at 9th-grade entry, on high school performance and the impact of these experiences on the transitions that students make from high school to adult roles. HSLS:09 was designed to help researchers, policymakers, and practitioners investigate the process of dropping out of high school and possible return to school or pursuit of alternative credentials; the school experience and academic performance of English language learners; the nature of the paths into and out of science, technology, engineering, and mathematics (STEM) curricula and occupations; and the educational and social experiences that affect these outcomes, decisions, and experiences. The second follow-up extended the focus of the study to emphasize the transition of the cohort to postsecondary education—both baccalaureate and subbaccalaureate—and the workforce, including access to higher education and choice of postsecondary institution. To that end, the PETS-SR collection reflects a shift in focus to higher education by collecting postsecondary transcripts and financial aid records. The longitudinal design of HSLS:09 is illustrated in figure 2.

**Figure 2. Longitudinal design for the HSLs:09 9th-grade cohort: 2009–2025**

SOURCE: U.S. Department of Education, National Center for Education Statistics, High School Longitudinal Study of 2009 (HSLs:09).

### 1.3.1 Base Year, First Follow-up, 2013 Update, and Second Follow-up

The HSLs:09 base-year data collection took place in the 2009–10 academic year with a randomly selected sample of fall-term 9th-graders in more than 900 public and private high schools with both 9th and 11th grades.<sup>3</sup> Students completed an in-person mathematics assessment focused on algebraic reasoning and a web-based survey that included items on educational experiences, sociodemographic background, educational expectations, and their perceptions of the value of science and mathematics as a subject area and as a vocation. Students’ parents, principals, and science and mathematics teachers, as well as their school’s lead counselors, completed surveys on the phone or on the Web.

The first follow-up of HSLs:09 took place in the spring of 2012, when most sample members were in 11th grade. The students were again assessed in mathematics, and

<sup>3</sup> Types of schools that were excluded from the sample based on the HSLs:09 eligibility definitions are described in the discussion of the target population in the *HSLs:09 Base-Year Data File Documentation* (see chapter 3, section 3.2.1) (Ingels et al. 2011).

they again completed a questionnaire. The first follow-up survey explored topics such as high schools attended, grade progression, school experiences, plans and preparations for the future transition out of high school, math and science aptitude and engagement, and extracurricular participation. Contextual data were collected from a subsample of parents and from school administrators and counselors. While re-administration of the counselor questionnaire occurred only in the base-year schools, administrator questionnaires were administered at base-year schools as well as the schools to which students had transferred.

The 2013 Update collection took place from June through January 2014. The 2013 Update was designed to collect information on the cohort's postsecondary plans and choices at the completion of high school (for most of the cohort). More specifically, information was collected about high school completion status, applications and acceptances to postsecondary institutions, education and work plans for the fall of 2013, financial aid applications and offers, choice of institution, and employment experiences. As part of the 2013 Update, high school transcripts were also collected in the 2013–14 academic year. Records matching (e.g., college admissions test scores, Free Application for Federal Student Aid data, GED<sup>4</sup> data) also contributed to the dataset.

The second follow-up data collection, conducted between March 2016 and January 2017, was designed to collect information from the cohort approximately 3 years after the modal high school completion date. At that time point, sample members may have been engaged in various activities, such as enrollment in postsecondary education, employment, serving in the military, volunteering, interning or getting other job-related training, and starting a family. Some sample members may have only recently received, or may still have been working toward, a high school credential. The survey explored a variety of topics that include, but are not limited to, high school completion and experiences, enrollment history and future enrollment plans, employment and unemployment history, family and home-life characteristics, and personal characteristics (e.g., disabilities, sexual orientation and gender identity, and civic engagement). The second follow-up survey also collected information on topics addressed in previous data collections, such as experiences, influences, and constraints on decision-making about postsecondary education, majors, and occupations with an emphasis on STEM fields.

---

<sup>4</sup> The GED credential is a high school equivalency credential earned by passing the GED test, which is administered by GED Testing Service. See <https://www.ged.com> for more information on the GED test and credential.

### 1.3.2 PETS-SR

**Postsecondary Education Transcript Study.** The HSLS:09 Postsecondary Education Transcript Study is the sixth Postsecondary Education Transcript Study (PETS) of high school cohorts. The first in the series (NLS:72) took place in 1984 and was followed by the HS&B sophomore cohort (1993), HS&B senior cohort (1986), NELS:88 (2000), and ELS:2002 (2013) PETS collections. Detailed PETS data files exist for NELS:88 and ELS:2002. Postsecondary education transcript studies have also been done in connection with the Beginning Postsecondary Students (BPS) and Baccalaureate and Beyond (B&B) longitudinal studies. A fundamental difference is that BPS and B&B are studies of students selected from a nationally representative sample of postsecondary institutions, while the high school studies are centered on a grade-cohort-based secondary school sample. In addition, BPS captures the full range of students entering postsecondary education for the first time, including students who begin their postsecondary education later in life, and the high school studies miss these late entrants if they begin their postsecondary education outside the study's time frame. Likewise, B&B is representative of baccalaureate recipients, and studies such as HSLS:09 and ELS:2002, which lack both late entrants and late completers, are not. Another key difference is that BPS and B&B also include postsecondary students who did not attend high school in the United States.

As an official institution record, the postsecondary transcript is a more reliable source of data regarding academic performance than is a student's self-report. The postsecondary transcript collection for HSLS:09, designed similarly to that conducted for ELS:2002 and BPS:04/09, provides much-needed information on the undergraduate experiences of 2009 ninth-graders who pursue postsecondary education in the years following high school. The combination of transcript data and other study data collected through interviews, by matching the sample to external data sources, and to student financial aid records collection allows researchers to analyze paths taken by cohort members as they begin undergraduate education. Postsecondary transcripts provide a wealth of data on enrollment, including degree or certificate program, terms enrolled, dual enrollment status, course intensity when enrolled, and fields of study. Furthermore, transcripts provide coursetaking details, including subjects taken and credits and grades earned. These data provide important links among the sample members' secondary academic performance, plans and expectations, and pathways into the workforce.

**Student Financial Aid Records Collection.** Previous secondary longitudinal studies collected student financial aid data from federal aid databases that detail federal aid exclusively. As such, a complete picture of all sources of student financial aid data, including both federal aid and nonfederal aid, has been lacking in the

secondary longitudinal studies, constituting a limitation in the utility of the study for analyses related to receipt of financial aid. Availability of financial aid is important at all points in the postsecondary process, including initial access and choice, persistence, transfer, and ultimate educational attainment. The financial aid data collected from the institutions attended by HSLS:09 sample members greatly increases the analytic utility of HSLS:09. Cumulative aid and debt can be calculated with scholarship, fellowship, grant, and loan amount data. The financial aid records collection also yields detailed information about students' enrollment patterns, degree or program of study and progress toward degree, and costs of attendance.

### **1.3.3 Research and Policy Issues and Analytic Levels**

This section broadly describes the research and policy questions targeted by HSLS:09, especially those related to transitions from high school to adult roles and experiences surrounding STEM. This section summarizes the conceptual model that was developed in the base year and that shaped many of the survey questions asked in previous collections. The primary focus of this section, however, is to describe the research and policy uses and analytic levels of the PETS-SR collection.

HSLS:09 is a general-purpose study, designed to serve multiple policy objectives rather than to test a specific hypothesis. The goal of HSLS:09 is better understanding of the relationship between earlier educational experiences, starting at 9th grade, and high school performance and the relationship of these experiences with the transitions that students make from high school to adult roles. HSLS:09 will help researchers and policy analysts investigate the features of effective high schools; growth in academic achievement, especially in mathematics;<sup>5</sup> factors related to dropping out of school and possible return to school or pursuit of alternative credentials; the school experience and academic performance of English language learners; the nature of the paths into and out of STEM curricula and occupations; and the educational and social experiences that affect these outcomes, decisions, and experiences.

The research agenda was guided by a conceptual model that was developed in the base year and shaped questionnaire content in both in-school rounds (i.e., fall 2009 base year and spring 2012 first follow-up). This model uses the student as the fundamental unit of analysis and attempts to identify factors that lead to academic goal setting and decision-making. It traces the many influences, including motivation, interests, perceived opportunities, barriers, and costs, on students' values and

---

<sup>5</sup> HSLS:09 includes an assessment of algebraic reasoning that measures achievement growth in the span between high school entry in the fall of 9th grade and the spring term of the junior year of high school for most cohort members (i.e., those in modal grade progression).

expectations that factor into their most basic education-related choices. Details of the conceptual model can be found in *HSLs:09 Base-Year Data File Documentation* (Ingels et al. 2011).

The 2013 Update and the second follow-up built on the information collected during the base-year and first follow-up collections. The 2013 Update collected information on the cohort's postsecondary plans and choices, gathered at, for most of the cohort, completion of high school. More specifically, information was elicited concerning high school completion status, applications and acceptances to postsecondary institutions, education and work plans for fall 2013, financial aid applications and offers, choice of institution, and employment experiences, as well as high school transcripts. The second follow-up was designed to collect information on the cohort's pursuit of postsecondary education, entry into the workforce, and family formation. Furthermore, the addition of postsecondary student financial aid records and postsecondary academic transcript information provides a continuous longitudinal record of courses taken, credit accrual, and grades through postsecondary years.

In total, the breadth of the study design supports researchers in exploring a multitude of analytic interests and policy issues. Several examples of lines of investigation are outlined below.

**Research and policy uses: base year and first follow-up.** Many topic areas can be investigated within the high school context. These areas include the process of dropping out or stopping out of high school (e.g., taking a temporary break), the resilience of students who persist despite multiple risk factors, the educational and occupational trajectories of students who remain in school but take extra time to graduate, achievement gains in mathematics and the correlates of academic growth, the role of family background and the home education support system in fostering students' educational success, the features of effective schools, and differential access to and engagement in various educational opportunities.

**Research and policy uses: 2013 Update.** The 2013 Update was administered in the last half of 2013. For students who graduated on time, the timing of the data collection corresponded to collection immediately after completion of secondary school. The 2013 Update questionnaire consisted of objective questions that could validly be completed either by parent or student; there was no preference for which respondent should complete the relatively brief survey. It was designed to elicit critical, time-sensitive data about how students and their parents make decisions about postsecondary choices. The 2013 Update provided information about status in summer/fall after the normative high school graduation, including educational status

(e.g., high school completion, continued high school enrollment, high school dropout status, and postsecondary attendance); work status; postsecondary education applications and financial aid; and work experiences.

**Research and policy uses: high school transcripts.** Data from the HSLS:09 High School Transcript component encompass coursetaking for grades 9–12. High school transcript data files can be analyzed on their own as stand-alone restricted-use files and can also be combined with the survey and assessment data for analysis.

High school transcript data from the secondary longitudinal studies also may be linked to postsecondary transcripts for high school cohort members who went on to postsecondary education or who were enrolled concurrently in postsecondary courses while in high school, known as “dual enrollment,” thus providing information for analyses relating academic preparation and experiences in high school to coursetaking and attainment in higher education (Adelman 2006). At the high school level, evidence from HS&B (Cool and Keith 1991; Meyer 1998), NELS:88 (Rock and Pollack 1995), ELS:2002 (Bozick and Ingels 2008), and the National Assessment of Educational Progress (NAEP) (Chaney, Burgdorf, and Atash 1997) suggests strong relationships between mathematics achievement and higher-level coursetaking.

**Research and policy uses: second follow-up survey.** Because most sample members in the 2016 second follow-up were 3 years beyond high school graduation, it is possible to study such topics as postsecondary education, entry into the workforce, and family formation.

The chief education-related foci of the second follow-up were access to postsecondary education, choice of postsecondary institution, and attainment of subbaccalaureate credentials. Early persistence and transfer from one postsecondary institution to another can also be studied. These topics of focus are asked of students who differ by postsecondary institution type and sector (e.g., public and private 2-year and 4-year institutions) attended; intensity of attendance (e.g., full-time versus part-time); whether enrollment was at the “first-choice” institution; and the institution’s location (e.g., urban, suburban, or rural; near home or distant). A student’s choice of postsecondary institution reflects institutional characteristics such as perceived academic quality or reputation, cost of attendance, and academic program offerings—all of which were captured in the 2013 Update and the second follow-up. The timing of the second follow-up also offered a window into attainment of 2-year degrees, postsecondary certificates, and certifications, whether granted by public institutions such as community colleges or by for-profit schools. The timing also provided an opportunity to view the transition from community

college settings to 4-year programs for those sample members whose pathway treats 2-year institutions as a stepping-stone to 4-year institutions. Other topics that can be explored include family formation; early occupational choice, with an emphasis on STEM fields; and labor market experiences.

**Research and policy uses: postsecondary education transcript study and student financial aid records collection.** In addition to information obtained from sample members who participated in the second follow-up survey, the PETS-SR collection, beginning in 2017, entailed collecting data from institutions and file matching to external sources. Financial aid data were collected from the institutions attended by HSLS:09 sample members, and federal student loan records were obtained from file matching. The financial aid records collected from the institutions attended by HSLS:09 sample members greatly increase the analytic utility of HSLS:09 and yields detailed information about students' enrollment patterns, degree or program of study, progress toward degree, and cost of attendance. The postsecondary transcript data cover postsecondary coursetaking through December 31, 2016, and provide detailed information on students' academic experience, including coursetaking, academic performance, credit accumulation, enrollment periods, and transfer between institutions.

As mentioned above, data from high school transcripts from secondary longitudinal studies may be linked to data from postsecondary transcripts for high school cohort members who enrolled in postsecondary education or who were dual enrolled in postsecondary courses while still in high school, thus providing information for analyses relating academic preparation in high school to coursetaking and attainment in postsecondary education (Adelman 2006).

**Research and policy uses: summary.** HSLS:09 helps researchers, educators, and policymakers understand outcomes associated with the 9th-grade cohort's continued academic, social, and interpersonal growth in and after high school. It illuminates the transitions from secondary and postsecondary education to the workforce. It also captures students' choices about access to and persistence in STEM courses and majors, or alternative (i.e., non-STEM) educational and career pathways. Finally, it helps identify and describe the characteristics of educational institutions and curricula that are related to student outcomes in adulthood, such as family formation (e.g., how prior experiences in and out of school relate to marital or parental status and how marital or parental status affects educational choice, persistence, and attainment); and characteristics of individual students associated with key outcomes, including how language-minority, low-socioeconomic status (SES), disability, racial/ethnic-minority, and at-risk status are associated with education and labor market outcomes for young adults.

**Analysis levels and design considerations.** The base-year HSLS:09 data can be analyzed cross-sectionally at both the student and the school levels. Fall 2009 entering high school freshmen can be descriptively profiled using the HSLS:09 nationally representative student sample. Analysis at the school level is also possible, supported by the HSLS:09 nationally representative sample of high schools with 9th and 11th grades.<sup>6</sup> HSLS:09 obtained information about the base-year schools from several sources: a school administrator questionnaire; school characteristics variables taken from the sampling frame consisting of the NCES Common Core of Data (CCD) and Private School Universe Survey (PSS); and the school's course offerings, as listed in school catalogs collected in the High School Transcript study.

In addition to the national samples of high schools and fall 2009 9th-graders, the data support analysis of 10 state-representative samples: California, Florida, Georgia, Michigan, North Carolina, Ohio, Pennsylvania, Tennessee, Texas, and Washington. The state samples pertain to the public sector only, and the national sample includes Catholic and other private schools.

Beyond the base year, HSLS:09 is representative of fall 2009 entering high school freshmen, who were followed up with 2 years after their 9th-grade year (first follow-up), the summer after the majority finished high school (2013 Update), and 3 years after the majority finished high school (second follow-up). HSLS:09 did not refresh the student sample; therefore, HSLS:09 is, for example, not representative of 11th-graders. Also, the representativeness of the school sample is lost after the base year.

HSLS:09 attempts to preserve the best design features of the predecessor high school longitudinal studies, while updating and improving upon those prior studies. The data collection points for HSLS:09 were chosen for their research value, considered independently of the data collection points employed in earlier secondary longitudinal studies. The base-year 9th-grade starting point was designed to capture—like NELS:88, which started in 8th grade—the transition into high school. It does so without the financial costs of following a sample in which 95 percent of the cohort had changed schools by the time of the first follow-up 2 years later, as experienced in NELS:88.

The HSLS:09 first follow-up took place when most students were in the spring term of 11th grade. It has often been observed that students in the spring of their senior year are disengaging from high school and not highly motivated to complete low-stakes assessments and questionnaires. Much thought has been given to improving students' participation and effort (e.g., as in NAEP, which traditionally has

---

<sup>6</sup> Researchers should note that, due to disclosure protections, relatively few school-level analyses can be done with the public-use files; for most purposes, the restricted-use files are required.

conducted 12th-grade as well as 4th- and 8th-grade assessments [see StandardsWork 2006]). One possible approach to addressing this problem is to move the testing point to spring of 11th grade, the strategy embraced by HSLS:09.

The timing of the 2013 Update—the last half of 2013 after (modal) graduation—also reflects a conscious choice. Earlier studies had data collections in the spring term—as early as January and February—of the senior year in high school, a time point at which many sample members had yet to make final decisions about postsecondary schooling or work. Much of the information about the decision process and its outcomes had to be collected, if at all, at the time of a follow-up 2 years after the senior year, when recollection of process details, including acceptances, rejections, and financial aid offers, had diminished. The Update’s timing strengthens the HSLS:09 longitudinal design by collecting decision information immediately following typical graduation.

The timing of the second follow-up, with student survey administration beginning in 2016, likewise was based on specific research considerations. In the past studies, the interval between high school graduation and the follow-up questionnaire was 2 years. For HSLS:09, the interval was 3 years. One benefit of this longer interval was the opportunity to obtain better information on postsecondary education persistence and subbaccalaureate attainment. A second benefit was that, at the time of the second follow-up, the subsets of HSLS:09 second follow-up students and BPS:12/14 first follow-up students who were immediate postsecondary entrants were aligned in terms of the amount of time that had elapsed since beginning postsecondary education. Both BPS:12/14 and HSLS:09 immediate postsecondary entrants were followed 3 years after first enrollment. Alignment of the two cohorts allowed for collection of postsecondary transcripts and student financial aid records to be conducted for both studies in tandem, thereby realizing efficiency gains.

Although HSLS:09 offers the design benefit of important new measurement points, a trade-off should be noted. Specific cross-cohort comparisons cannot be made with the earlier secondary longitudinal studies. Nor can comparisons be made with the high school transcript studies of NAEP. HSLS:09 is based solely on a fall 9th-grade cohort, whereas the prior longitudinal studies were based on spring-term 8th-, 10th-, or 12th-grade cohorts (see figure 1). NAEP transcripts were collected only for graduating seniors and are nationally representative for that population. Similarly, the links between NAEP, NELS:88, and ELS:2002 mathematics assessments cannot be replicated within the HSLS:09 design.

A final point about the comparative structures of HSLS:09 and its two most recent predecessor studies pertains to sample “freshening,” a device for cost-efficiently

generating multiple grade-representative cohorts during a longitudinal study. As mentioned above, HSLS:09 includes only a single cohort, not two (grades 10 and 12 as in ELS:2002) or three (grades 8, 10, and 12 as in NELS:88); the 9th-grade student sample is the sole cohort across all rounds. The earlier studies freshened the sample to represent later grades. This was done for a compelling reason: to facilitate cross-cohort comparisons (e.g., trends among high school seniors in 1972, 1980, and 1992). Because HSLS:09 has no specific cross-cohort comparison points within the family of NCES secondary longitudinal studies, the traditional rationale for freshening does not apply. Freshening was also problematic because the 9th-grade sample does not represent all, or nearly all, 9th-graders—schools were eligible if and only if they had both a 9th grade and an 11th grade at the time of sampling.

## Chapter 2. Sample Design

This chapter provides details of the sample design employed for the High School Longitudinal Study of 2009 (HSLs:09). Sections 2.1, 2.2, 2.3, and 2.4 summarize the school and student sampling used for the base year, first follow-up, 2013 Update and High School Transcript study, and second follow-up, respectively. The student sample for the PETS-SR collection is described in section 2.5.

### 2.1 Base-year Sample Design

**Selection of the school sample.** HSLs:09 employed a stratified, two-stage random sample design with primary sampling units defined as schools selected in the first stage and students randomly selected from the sampled schools in the second stage. The HSLs:09 target population of schools was defined in the base year as regular public schools, including public charter schools, and private schools in the 50 states and the District of Columbia that provided instruction to students in both the 9th and 11th grades as of fall 2009.<sup>7</sup> For details of the rules for school inclusion or exclusion, see the *HSLs:09 Base-Year Data File Documentation* (Ingels et al. 2011). A total of 944 of 1,889 eligible schools participated in the base year, resulting in a 55.5 percent base weighted school response rate.

Although HSLs:09 was designed to be representative of 9th-grade students in the 2009–10 school year in study-eligible schools across the United States, it also supports construction of select state-level estimates for students enrolled in 9th grade in public schools in the fall of 2009. In particular, in response to a request from the National Science Foundation for representative estimates within certain states, the design was augmented with additional sample schools to support the revised study objectives within 10 states: California, Florida, Georgia, Michigan, North Carolina, Ohio, Pennsylvania, Tennessee, Texas, and Washington. Additional information on construction of the HSLs:09 base-year school sample may be found in the *HSLs:09 Base-Year Data File Documentation* (Ingels et al. 2011).

---

<sup>7</sup> Note that some schools that meet these definitional criteria would not meet current federal legal definitions of high schools which must also offer grade 12 classes and grant diplomas, as per 20 USC § 7801 (28). See <https://www.govinfo.gov/content/pkg/USCODE-2018-title20/html/USCODE-2018-title20-chap70.htm>.

**Selection of the student and contextual samples.** The student target population contained all 9th-grade students as of fall 2009 who attended either regular public or private schools<sup>8</sup> in the 50 states and the District of Columbia that provided instruction in both 9th and 11th grades. This population is referred to as the “9th-grade cohort” in the subsequent discussions.

A sample of 26,305 students was randomly selected from the 944 participating schools in the base year. During base-year recruitment, 1,099 students (4.2 percent unweighted) were classified as study ineligible and excluded from the data collection rosters, yielding 25,206 study-eligible students. Student participants completed an in-school survey (85.7 percent base weighted) and mathematics assessment (83 percent base weighted).

Contextual information was collected on the student sample to describe the home and school environments. Home-life and background information was obtained through questionnaires completed by students’ parents. Administrator and counselor questionnaires provided school information. Teacher questionnaires, completed by science and mathematics teachers linked to the sampled student, captured information on teacher background and preparation, school climate, and subject-specific and classroom practices.

For additional information on selection of the HSLS:09 base-year student and contextual samples, please refer to the *HSLS:09 Base-Year Data File Documentation* (Ingels et al. 2011).

## 2.2 First Follow-up Sample Design

The first follow-up student target population is the same as defined for the base year.

**First follow-up student and contextual samples.** All 25,206 base-year study-eligible students—regardless of their response and enrollment status—were included in the first follow-up sample. Unlike prior NCES high school longitudinal studies, NELS:88 and ELS:2002, the HSLS:09 student sample was not freshened to include a representative later-grade cohort, such as 11th-graders in HSLS:09. Therefore, first follow-up estimates from the sample are associated only with the 9th-grade cohort

---

<sup>8</sup> The term “regular” refers to the setting and mode of instruction. Some examples of schools not considered regular are those that offer instruction in juvenile detention centers, schools that instruct only special education students, and schools where all the students may be homeschooled or where a mix of instructional modes is used (e.g., some students are homeschooled, some receive remote instruction, and some are in a common physical location).

2 1/2 years later, and not the universe of students attending the 11th grade in the spring of 2012.

Some students were deceased as of the first follow-up, withdrew from HSLS:09 prior to the first follow-up, or were determined to be study ineligible for HSLS:09 as of the first follow-up; 25,184 remained eligible as of the first follow-up.

The student questionnaire explored a variety of topics that include, but are not limited to, high school attendance, grade progression, school experiences, demographics and family background, completion of admission tests, college choice and characteristics, and high school coursetaking. Contextual information was collected for the student sample to describe their home and school environments. Home-life and background information was obtained through questionnaires completed by students' parents. The first follow-up parent questionnaires were administered to the parents of a random 48 percent subsample of students, whereas parent questionnaires were sought for all students in the base year. School information was obtained through the administrator and counselor questionnaires; however, administrator data were collected at both the base-year schools and the schools to which sample members transferred. Counselor data were collected in the first follow-up only from base-year high schools. For additional information on selection of the HSLS:09 first follow-up student and contextual samples, see the *HSLS:09 Base Year to First Follow-Up Data File Documentation* (Ingels et al. 2013).

## 2.3 2013 Update and High School Transcript Study Sample Design

In the 2013 Update, students or their parents responded to a survey in which information was collected on the student sample to describe the student's high school completion status, postsecondary education and work plans, college application experiences, and work experiences. In addition, school personnel in base-year schools and other schools identified during data collection supplied high school transcripts for HSLS:09 students from all schools that these students had attended.

As of the 2013 Update, 25,168 remained eligible and 25,167 remained eligible as of the High School Transcript data collection.

Of the first follow-up eligible and fielded sample members, 1,767 were not fielded for the 2013 Update. The majority of these sample members were nonrespondents in both the base year and the first follow-up. Additionally, some sample members were not fielded for the 2013 Update because they withdrew from the study. Information

on selection of the HSLS:09 2013 Update sample appears in the *HSLS:09 2013 Update and High School Transcript Data File Documentation* (Ingels et al. 2015).

## 2.4 Second Follow-up Sample Design

The second follow-up fielded sample included 23,316 of the 23,401 sample members fielded and found eligible for the 2013 Update. The 85 sample members not fielded withdrew from the study between the end of the 2013 Update collection and the beginning of the second follow-up data collection or were found to be deceased.

## 2.5 Postsecondary Education Transcript and Student Records Sample Design

Among the 3,491 institutions reported to have been attended by sample members, it was determined that 220 institutions were ineligible because the institution had closed, because a sample member had reported a school that was not a postsecondary institution, or because all of the sampled students were reported as having not attended the institution. Hence, transcripts and student records were requested from 3,271 postsecondary institutions.

Sample members eligible for PETS-SR consisted of only those who were ever enrolled at an Integrated Postsecondary Education Data System (IPEDS)-participating postsecondary institution as of June 30, 2017. Prior to the start of data collection, there were 17,201 students identified as being enrolled based on responses to the 2013 Update and second follow-up surveys as well as National Student Loan Data System data (NSLDS) matching. These 17,201 students were fielded for PETS and SR data collection. During the course of data collection, it was determined that 328 of these cases were not eligible (i.e., did not attend, based on reports from the institution), resulting in 16,873 eligible fielded cases. After data collection, a match to National Student Clearinghouse (NSC) identified 474 additional eligible cases. However, 9 of these cases were found to be deceased. Therefore, for weighting purposes, 17,338 cases were eligible for the PETS component. Eligibility for the SR component is a subset of PETS; in particular, sample members eligible for SR were those who were enrolled in a postsecondary institution after the completion of high school or high school equivalency. Hence, by contraposition, sample members eligible for PETS but not SR were sample members who were at one time enrolled at an IPEDS-participating institution, but who were not enrolled following the completion of high school or a high school equivalency program. For weighting purposes, 17,230 cases were eligible for the SR component.

## Chapter 3. Data Collection Methodology and Results

### 3.1 Postsecondary Education Transcripts and Student Records Systems and Processes

#### 3.1.1 *Postsecondary Data Portal Website*

This section provides information about the Postsecondary Data Portal (PDP), which listed information about the study and served as a secure platform for institutions to upload requested electronic data (i.e., student records and transcripts).

The PDP website contained information about NCES sample surveys that collect data through the system, including research topics, the ways in which data would be used, answers to frequently asked questions, and confidentiality assurances. Contact information for the data collection Help Desk, project staff at RTI, and NCES project officers, as well as a link to the main NCES website, were also included on the website. From the credentialed-access portion of the website, authorized personnel from institutions could view the list of their sampled students, view detailed instructions for entering or uploading data, and enter or upload data.

Various security measures were incorporated into the website application to ensure strict adherence to NCES confidentiality guidelines, including

- a Secure Sockets Layer Certificate that ensured secure data transmission over the Internet;
- password protection of all data-entry modules;
- automated user log-out after 20 minutes of inactivity; and
- a secure file transfer mechanism, in which files uploaded to the secure website were immediately moved to a secure project folder accessible only to a subset of project staff.

### 3.1.2 *Institution Contacting Staff Training*

Institution contacting staff consisted of institution contactors (ICs) and quality control supervisors (QCSs) who were responsible for staff supervision. Prior to the start of data collection, these staff were trained over a 4-day period on the study's background, gaining cooperation, problem resolution, and collection and receipt systems. Training also included provision of answers to frequently asked questions and a review of confidentiality regulations.

### 3.1.3 *Institution Contacting and Recruitment*

Institutions attended by sample members were asked, in the same requests, to participate in the PETS and SR collections. If separate staff members were identified to provide different types of data, contact materials were directed to those staff as needed. Follow-up contacts occurred after the initial mailing to confirm receipt of the package and answer any questions about the study, as applicable.

Transcripts and student records were requested from 3,271 eligible, fielded postsecondary institutions. In addition, if an institution had copies of transcripts received from any transfer schools attended by the sample members, the transfer transcripts were requested as well. The requests included 29,766 transcripts and 29,633 student records covering 17,201 students.

Transcripts and student records were first requested in March 2017. Student records data collection ended in February 2018, and transcript collection ended in November 2018.

**Contacting institutions.** A web-based control system—the Institution Contacting System (ICS)—supported each step of the transcript and student records collections, including project management, communications, and tracking. The ICS was used to store and access data on students and track efforts to obtain their transcripts and student records data. Prior to the start of collection, the ICS was loaded with the institution sample including, when known, contact information for the institution. When needed, Internet searches were conducted to identify the director of the institutional research office and the registrar. Phone calls were then made to institutions to confirm or obtain appropriate contact information. At the start of data collection, a request packet was then sent to the institutional research director. In the absence of an office of institutional research, packets were sent to the chief administrator's office. Follow-up calls by trained ICs were placed about 2 days after the initial mailing to confirm receipt of the packet and to answer any questions about the study. Prompting calls were made and reminder e-mails sent, as needed,

throughout data collection. During initial conversations, institution staff members were asked to identify a coordinator to serve as the primary point of contact for the data collection as well as additional contact people for transcript collection and/or student records collection.

ICs also served on the Help Desk team to assist institution staff who called in or e-mailed questions. Incoming calls from institution staff were most often related to two areas: requests for an extension to the data submission deadline and assistance with the website such as with password resets and data uploads.

In addition to the data requests made to individual institutions, institution systems or groups of institutions were identified where data were supplied by one office or individual for all of the institutions. Contacts at the system level or corporate level with the access or ability to coordinate provision of this information vary, depending on the type of institution and the culture of the institution. Hence, data release may have been provided by or coordinated by a system-wide office of institutional research, the office of government relations, the student financial aid office, or an information technology group that is routinely charged with handling data requests for individual institutions. This strategy of utilizing centralized contacts has been successful in increasing the efficiency of this and previous institution data collections and minimizes burden by removing the need to contact each institution within the system separately. It was often the individual institutions themselves that pointed us to the appropriate system-wide contact.

**Mailings.** As mentioned above, institution staff received a single letter informing them of the request for transcripts and student records. They also received a transcript request packet and a student records packet that could be passed along to another individual at the institution, if needed. See appendix C for examples of data collection notification materials.

Request packets for the transcript collection and the student records collection included:

- letters introducing the study, requesting data, and providing information regarding how to log into the study's secure website;
- instructions for providing data; and
- a brochure.

Additional information, such as data confidentiality and the study's compliance with the Family Educational Rights and Privacy Act, was available on the PDP.

**Submission modes.** Institutions were provided with multiple options for how they could submit the data requested.

With respect to submission of transcript data, modes included (1) upload, by institution staff, of electronic transcripts for sampled students to the secure PDP website; (2) delivery of electronic transcripts via secure File Transfer Protocol (FTP); (3) submission of electronic transcripts as encrypted attachments via e-mail; (4) request/collection by RTI of electronic transcripts via a dedicated server at the NSC for institutions that already use this method; (5) submission of electronic transcripts via eSCRIP-SAFE, in which institutions send data to the eSCRIP-SAFE server by secure internet connection after which they can be downloaded only by a designated user; (6) transmission of transcripts via a secure electronic fax after a test submission of nonsensitive data confirms that the institution has the correct fax number; and as a last resort, (7) delivery of redacted transcripts via FedEx.

In addition to transcripts, other information from each institution was requested to facilitate transcript keying and coding. Institutions were asked to provide academic calendar and grading system information. Course catalogs were also sought as reference material to code courses. The majority of the catalogs were obtained from institution websites.

For student records, three options were offered to institutions for providing the data, similar to those used for the 2015–16 National Postsecondary Student Aid Study (NPSAS:16). Institution coordinators were invited to select the delivery method that was most convenient for the institution.

The options for providing student records data included (1) institution staff keying data into the PDP's web-based data-entry interface by student, by year; (2) institution staff keying data into an Excel workbook that is preloaded with student identifying information and then uploading it to the PDP; and (3) institution staff creating CSV (comma-separated values) data files according to study specifications and uploading them to the PDP.

**Quality control and follow-up.** During the collection period, data collection staff members met weekly to review progress, ask questions, and discuss any issues. Project staff used daily monitoring reports to review potential errors in received data, and ICs recontacted institutions to resolve issues or request additional or replacement data.

## 3.2 Postsecondary Education Transcripts-Specific Processes and Quality Control

### 3.2.1 Data Receipt Procedures

The initial transcript check-in procedure was designed to log the receipt of materials into the Data Receipt System (DRS) as they were received each day. Transcripts and supplementary materials received from institutions (including course catalogs) were inventoried, assigned unique identifiers based on the IPEDS ID, reviewed for problems, and logged in the DRS. Received transcripts were reviewed by project staff for completeness. ICs then contacted the institutions to prompt for missing data and to resolve any problems or inconsistencies.

Transcripts received in hardcopy form were subjected to a brief review prior to recording their receipt. Receipt control clerks checked transcripts for completeness and reviewed transmittal documents to ensure that transcripts were received for each of the specified sample members. The disposition code for transcripts received was entered into the DRS. Course catalogs were also reviewed and their disposition status updated in the system for cases in which this information was necessary and not available through institution websites. Hardcopy course catalogs were sorted and stored in a secure facility at RTI, organized by institution. The procedures for electronic transcripts were similar to those for hardcopy documents—receipt control personnel, assisted by programming staff, verified that the transcript was received for the requested sample member, recorded the information in the receipt control system, and verified that a readable, complete electronic transcript was received.

Data-processing staff were responsible for (1) associating files with the sending institution; (2) associating files with the correct sampled students, at which point the transcript file was given an ID number; (3) reviewing the transcript files to identify missing, incomplete, or indecipherable transcripts; and (4) assigning appropriate problem codes for missing and problematic transcripts as well as providing detailed notes regarding each problem to facilitate follow-up by ICs and other project staff. Project staff used daily monitoring reports to review the transcript problems and to identify approaches to resolve the issues. Web-based collection allowed timely quality control, as RTI staff were able to monitor data quality for participating institutions closely and on a regular basis. When institutions called the Help Desk for technical or substantive support, the institution's data could be queried directly to aid communication with the institution and effectively resolve problems. Transcripts were shredded or destroyed after the transcripts were keyed, coded, and quality checked.

### 3.2.2 Transcript Keying/Coding System and Keyer/Coders

This section describes the transcript Keying and Coding System (KCS), which facilitated the efficient and secure capture of data from student transcripts. The section also provides an overview of the training provided to transcript keyer/coders (KCs).

Once received, transcripts were keyed and coded using the KCS, a web-based platform for data entry that facilitated the efficient and secure capture of data from student transcripts. The application included five main pages in which student-level data were stored: Case Information, Schools and Terms, Tests, Degrees, and Courses. For each page, project staff used the transcript and institution-specific course catalog to encode relevant data. If a datum was not present on the transcript, such as a test score, the field was left blank in the KCS.

1. *Case Information.* The Case Information page captured the student's name, address, date of birth, Social Security number, and high school graduation date. Depending on the transcript, project staff entered complete or partial information for each of these elements (e.g., last four digits of the Social Security number).
2. *Schools and Terms.* On the Schools and Terms page, HSLS:09 project staff confirmed that all schools appearing on a transcript were captured in the KCS, including the school issuing the transcript and any transfer schools. Schools were preloaded for each student based on previously collected data—the 2013 Update and second follow-up surveys as well as NSLDS data matching—and identification during the receipt process. Staff members entered all academic terms in which the student was enrolled in at least one course and also entered the student's cumulative grade point average (GPA) on this page.
3. *Tests.* For all tests that appeared on the transcript, such as the SAT, the test name and score were captured on the Tests page.
4. *Degrees.* Any degree programs attempted or earned were entered on the Degrees page. If the degree was awarded, the date of receipt as well as any graduation honors were keyed. For each degree program, the field of study was keyed and coded on this page, including majors, minors, and concentrations.
5. *Courses.* HSLS:09 project staff entered course-specific information and coded course content on the Courses page. Elements included the term in which

the course was taken, course number, course name, grade earned, and credit or clock hours earned or attempted. Course attributes, such as lab, and noncourse credit, such as credit for an Advanced Placement (AP) exam, were keyed on this page.

Figures 3 and 4 show screenshots of the Degrees and Courses pages from the KCS, respectively. Once keyed, this information was presented in tabular form.

**Figure 3. Keying and Coding System Degrees page: 2018**

SOURCE: U.S. Department of Education, National Center for Education Statistics, High School Longitudinal Study of 2009 Postsecondary Education Transcript Study and Student Financial Aid Records Collection.

**Figure 4. Keying and Coding System Courses page: 2018**

SOURCE: U.S. Department of Education, National Center for Education Statistics, High School Longitudinal Study of 2009 Postsecondary Education Transcript Study and Student Financial Aid Records Collection.

To enhance the utility of the KCS and aid in validating the accuracy of captured data, the KCS included specific features, such as direct links to the transcript PDF and course catalog files, validated fields for value range and value types, and requirements that all fields of study and courses received a code.

### 3.2.2.1 *Transcript Keyer/Coder Training*

Keying and coding training took place in April 2017 and consisted of 12 KCs, two quality experts, and one project supervisor. Additional trainings were conducted in August and September 2017 to increase the number of KC staff. The project supervisor was the administrative manager of the KCs housed at RTP's Research Operations Center (ROC), a separate facility for coding and telephone interviewing staff. Quality experts were responsible for assisting with quality control during data collection, answering questions from KCs, and conveying more difficult questions to project staff as needed. Prior to the training, KCs signed confidentiality agreements and notarized affidavits and initiated Electronic Questionnaires for Investigations Processing (e-QIP) applications.

The training presented an overview of HSLS:09, including the role of KCs within the PETS study component. KCs were introduced to transcript and course catalog formatting and the data elements captured within the KCS. The functionality of the system was demonstrated followed by practice sessions. The training provided examples of common and difficult situations as well as how to capture the data properly in the KCS. The final day of training consisted of supervised practice and an exam wherein KCs keyed and coded an entire transcript on their own.

Quality Circle meetings were held weekly—and later, biweekly—to inform the KCs of updated protocols and to discuss common issues confronted by the KCs. These meetings provided an opportunity for KCs to ask general or specific keying and coding questions to project staff. Using information from these meetings, guidance documents were updated to reflect new best practice guidelines or updates to existing procedures. All training and guidance documents were made available electronically to KCs for the duration of the project.

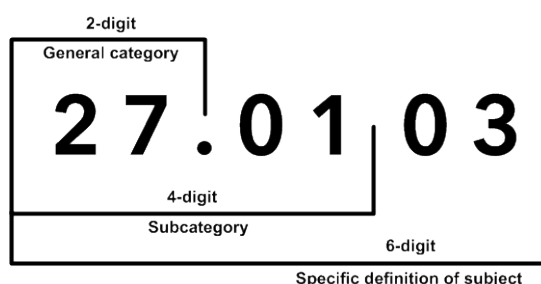
### 3.2.3 *Coding Taxonomies*

To standardize institution names, fields of study, and course content collected from transcripts, the KCS used three coding taxonomies. Postsecondary institutions at which students were enrolled were coded using IPEDS, developed by NCES (<https://nces.ed.gov/IPEDS/>). Students' fields of study data, captured as majors, minors, and concentrations, were coded using the 2010 Classification of Instructional Programs (CIP) taxonomy developed by NCES (<https://nces.ed.gov/ipeds/cipcode/>). Course content was coded using the 2010 College Course Map (CCM) (Bryan and Simone 2012). The CCM is a taxonomy for coding postsecondary education courses and it extends the list of codes contained in the CIP. The CIP contains codes for instructional programs, whereas the CCM also

contains course-specific codes. For example, the CCM contains a code for abstract algebra because this is a course rather than a program of study, such as computational mathematics. In total, the CCM contains 463 course codes not included in the CIP.

The CIP and CCM share a six-digit code structure in which the first two digits identify the general category, the first four digits indicate the subcategory within the general category, and the six-digit code provides the specific definition of the field of study or course content. Figure 5 presents a visual representation of the structure of the codes.

**Figure 5. Code diagram: 2018**



SOURCE: U.S. Department of Education, National Center for Education Statistics, High School Longitudinal Study of 2009 Postsecondary Education Transcript Study and Student Financial Aid Records Collection.

This structure includes codes that can be described as general, specific, or other. Eighteen percent of the codes in the CCM are general codes; these typically end with 00 or 01 and represent undifferentiated courses within the main category or subcategory. For example, 27.0101 is “Mathematics, General.” More specific subjects are represented with codes ascending from the general code, such as 27.0103, “Analysis and Functional Analysis.” In the aforementioned example, 27 represents the general category of “Mathematics and Statistics,” the 4-digit subcategory 27.01 encompasses a number of pure mathematics courses, while the full 6-digit code specifically defines “Analysis and Functional Analysis.” Eleven percent of the codes are described as other and represent those subjects that are not general and not covered in a specific code. These codes end with 99; for example, 27.0399, “Applied Mathematics, Other.”

### 3.2.4 Transcript Keying/Coding Quality Control

The transcript keying/coding task included multiple quality-control procedures, such as double-coding with arbitration, double-keying, as well as reviewing and upcoding of unknown institutions and “other, specify” fields.

### 3.2.4.1 Double-Coding with Arbitration

Double-coding was performed within the KCS on 10 percent of courses from each institution. In cases where students took fewer than 10 courses at an institution, all of the courses were double-coded. To complete double-coding, a second KC reviewed the course data from the transcript and selected a code. The second KC did not have knowledge of the code chosen by the first KC.

The results of double-coding were used to evaluate the reliability of course coding and to offer feedback to improve course coding. For HSL:09, 30,039 courses were double-coded and the results of agreement between KCs measured the inter-rater reliability rate for course coding.

Cohen's kappa statistic was used to assess inter-rater reliability between the two KCs at the three levels of the CCM code: the two-digit general category, the four-digit subcategory, and the specific six-digit course code. Measuring the proportion of agreement between raters, above what would be expected by random chance, a kappa score of 0.81–1.00 is considered “almost perfect agreement,” 0.61–0.80 is “substantial agreement,” and 0.41–0.60 is “moderate agreement” (Cohen, 1960). The kappa value for two-digit agreement between coders was 0.83, indicating almost perfect agreement on general course categories. The kappa value for four- and six-digit agreement between coders was 0.73 and 0.61, respectively, indicating substantial agreement at the subcategory and specific code levels.

Arbitration was conducted for those courses in which the first and second KCs did not select the same code. The arbiter, a member of the project staff with extensive knowledge of the taxonomy and coding guidance, reviewed the course information and selected a course code. The arbiter had access to the codes selected by both KCs and could choose a code that agreed with either KC or a third code that differed from both KCs. Results were used to provide feedback to KCs and to develop additional guidance on selecting the best-fitting codes.

### 3.2.4.2 Double-Keying

Double-keying was performed on a random 10 percent sample of the transcripts completed by each KC. A subset of transcript items was double-keyed; this subset included degree programs, terms, demographic information, and tests. As with double-coding, the results were assessed using Cohen's kappa statistic to evaluate the reliability of keyed data. The kappa value for degree programs was 0.95 (almost perfect agreement), for terms it was 0.99 (almost perfect agreement), for demographic information it was 0.81 (almost perfect agreement), and for tests it was 0.45 (moderate agreement). The lower agreement rate for tests is partly due to

scoring changes. Standardized tests, such as the ACT, have changed the scoring rubric over time, and KCs had to choose the scoring rubric that matched the transcript using a drop-down menu that included, for example, both the current ACT rubric and old scoring rubric. Institutions often reported standardized test scores in highly variable ways, and in many cases did not present subscores (e.g., Math, Writing) separately, making it difficult to identify which scoring rubric was used. In addition, although the components of the ACT have not changed, the range for the writing score changed in 2016. If a transcript did not contain a writing score, it is possible that the KC chose the incorrect option to enter the test score(s).

### **3.2.4.3 Review and Upcoding of Unknown Institutions**

Courses which appeared on transcripts were associated with the institution at which the courses were taken. Therefore, when courses were taken at an institution other than the one which sent the transcript, the other institution was coded using an IPEDS code. Project staff were unable to code some institutions due to inadequate or unclear data from the transcript, or due to the institution being located internationally and thus not included in IPEDS. Institutions that were not coded were reviewed by project staff to determine if ancillary information could be used to assign an IPEDS code. For example, because HSLS:09 attempted to collect transcripts from all institutions the student attended, additional transcripts often clarified the previously unknown institution. After final reconciliation, six institutions remained uncodable. Uncodable institutions are indicated by an IPEDS ID that begins in 8.

### **3.2.4.4 Review and Upcoding of Variables with “Other, specify”**

A number of variables within the KCS contained an “other, specify” option. The text strings entered for these options were reviewed and upcoded, as applicable. The KCS variables with an “other, specify” option included the following:

- tests (e.g., SAT);
- degree programs (e.g., diplomas and certificates);
- degree types (e.g., Bachelor of Education);
- graduation honors (e.g., magna cum laude);
- grades (e.g., E);
- term honors (e.g., dean’s list); and
- noncourse credits awarded (e.g., course credit for AP tests).

The text strings entered in these fields were upcoded by placing the response into an existing option or into a newly formed option when similar strings were identified. In instances where the string could not be upcoded, the value was left unchanged.

Table 1 shows the results of “other, specify” upcoding. The total number of cases is shown for each data element along with the number and percent that were upcoded.

**Table 1. Upcoding of “other, specify” responses: 2018**

<b>Data elements with “other, specify” option</b>	<b>Number of “other, specify” cases</b>	<b>Number upcoded</b>	<b>Percent upcoded</b>
Degree programs	195	170	87.2
Degree types	559	368	65.8
Graduation honors	179	109	60.9
Grades	21,554	17,623	81.8
Term honors	2,203	706	32.0
Noncourse credits awarded	3,815	3,452	90.5

SOURCE: U.S. Department of Education, National Center for Education Statistics, High School Longitudinal Study of 2009 Postsecondary Education Transcript Study and Student Financial Aid Records Collection.

Degree programs, such as Associate’s, Bachelor’s, and undergraduate certificates were also captured in the KCS. Almost 200 were categorized as “other,” of which 87 percent were upcoded to existing categories. Degree types included specific degrees such as Bachelor of Science. Of those that were categorized as “other,” 66 percent could be upcoded into existing categories. The KCS captured the grade received for each course, the majority of which were common letter grades (e.g., A-F, I, and W) and numeric grades. Uncommon grades were also captured, of which 82 percent were found to be equivalent to common letter or numeric grades. For example, a transcript may indicate a grade of “U,” which project staff later determine to be equivalent to “Withdrawal” after inspecting the institution’s course catalog. Note that not all institutions provide transcript keys which indicate the definition of nonnumeric or nonstandard grades. In addition to courses, the KCS captured credits earned for noncoursework activities, such as taking a test (e.g., AP or CLEP) or work or military experience. Among these credits identified as “other,” 90 percent were upcoded to existing categories. Graduation honors, such as cum laude, were also captured in the KCS. Among these, 179 were identified as “other” of which 61 percent were upcoded to existing categories. The KCS captured term honors (e.g., dean’s list, president’s list) as well, including nearly 2,200 term honors which required upcoding. Of these, 32 percent were successfully upcoded into existing categories.

Although reviewed, test data were not upcoded due to the high degree of variability in responses. Test data found on transcripts were more frequently entered with the “other, specify” option than with response options for known tests such as the SAT or ACT.

### 3.2.4.5 Review of “Other” Courses

As noted in section 3.2.3, the CCM includes courses designated as “other.” These course codes were used to code courses with descriptions that did not match specific course codes. Project staff reviewed the courses coded as “other” to determine if there were common subjects among the codes that would merit introduction of new codes to the taxonomy. Of all CCM codes used at least once, the median frequency was 49. This median was used as the threshold for adding a new code: if 49 instances of courses with the same subject could be identified within those coded with a particular “other” code, a new code would be added. However, a review of “other” codes did not identify any subjects that met this threshold; hence no new codes were introduced.

## 3.3 Postsecondary Education Transcripts Data Collection Results

This section provides the results of the transcript data collection, including the number of transcripts, number of keyed/coded transcripts, and number of courses coded, including those identified as uncodable.

**Institution-level transcript collection.** Table 2 provides institution participation rates by institution type. The fielded institution sample for the transcript collection included 3,491 distinct IPEDS institutions. As noted in section 2.5, of the 3,491 institutions, it was determined that 220 cases were ineligible because the institution had closed, because a sample member had reported a school that was not a postsecondary institution, or because all of the sampled students were reported as having not attended the institution. Of the remaining 3,271 institutions, 2,517 (77 percent) provided information. Across the institution types represented, participation in the transcript collection ranged from 97 percent at public 4-year, doctorate-granting institutions to 40 percent at private for-profit 2-year institutions. Some common reasons cited by institutions for not participating in PETS included lack of available staff to handle the request for transcripts and the timing of the transcript request.

**Table 2 Eligible institution participation, by institution type: 2018**

Number	Percent	Institution-level participation <sup>1</sup>	
		Number	Percent
<b>Total</b>	<b>3,271</b>	<b>2,517</b>	<b>76.9</b>
Institution type			
Public			
Less-than-2-year	55	25	45.5
2-year	787	719	91.4
4-year, non-doctorate-granting	304	266	87.5
4-year, doctorate-granting	334	323	96.7
Private nonprofit			
Less-than-4-year	63	29	46.0
4-year, non-doctorate-granting	514	458	89.1
4-year, doctorate-granting	369	344	93.2
Private for-profit			
Less-than-2-year	311	144	46.3
2-year	258	104	40.3
4-year	216	102	47.2

<sup>1</sup> An institution was considered a participant if it provided information for at least one student. A small number of the participating institutions are not represented in the institution type rows due to unknown institution sector.

SOURCE: U.S. Department of Education, National Center for Education Statistics, High School Longitudinal Study of 2009 Postsecondary Education Transcript Study and Student Financial Aid Records Collection.

**Student-level transcript collection.** At the student-level, postsecondary transcripts were pursued for 17,201 sample members, however 328 sample members were found to be ineligible, leaving 16,873 sample members who were eligible or had unknown eligibility. A transcript was received from at least one institution for 13,160 sample members (78 percent). Table 3 shows the transcript collection results at the student level.

**Table 3. Student-level transcript collection results: 2018**

Student sample	Number	Percent
<b>Total fielded sample</b>	<b>17,201</b>	<b>†</b>
Fielded, eligible or unknown eligibility	16,873	98.1
Eligible, at least one transcript received <sup>1, 2</sup>	13,160	78.0
Eligible, transcript nonrespondents <sup>2</sup>	3,280	19.4
Unknown eligibility	433	2.6
Fielded, ineligible	328	1.9

† Not applicable.

<sup>1</sup> A student was considered a transcript respondent if a transcript was received from one or more institutions and the transcript contained a course, term, or a degree program.

<sup>2</sup> Of the total sample ( $n = 25,206$ ), 898 sample members were known to be eligible for the transcript collection but were not fielded because they were not fielded in the second follow-up.

NOTE: Detail may not sum to totals because of rounding.

SOURCE: U.S. Department of Education, National Center for Education Statistics, High School Longitudinal Study of 2009 Postsecondary Education Transcript Study and Student Financial Aid Records Collection.

## 3.4 Student Records-Specific Processes and Quality Control

### 3.4.1 Student Records Instrument

The HSLS:09 student records instrument collected data about sample members' postsecondary education for up to six academic years, 2011–12 through 2016–17.

The instrument content was divided into five sections:

1. Institution Information, which collected each institution's enrollment terms during each academic year.
2. General Student Information, which collected students' demographic characteristics.
3. Enrollment, which collected the information about students' degree program, major(s), class level, and enrollment intensity at the institution for each academic year.
4. Budget, which collected the budgeted cost of attending the institution for each academic year.
5. Financial Aid, which collected all financial aid awarded to the student for each academic year. This section included federal, state, institution, graduate, and any private or other government awards.

For the full list of items collected in each section, see appendix B.

HSLs:09 is one of several recent studies that have collected student records information through the PDP; these include NPSAS:16, the 2011–12 Beginning Postsecondary Students Longitudinal Study (BPS:12), and the 2017–18 National Postsecondary Student Aid Study Administrative Collection (NPSAS:18-AC). Student records items largely remained stable across these NCES studies to ensure familiarity and consistency for each collection, as well as minimize burden on the institutions because they can use prior-round procedures for providing the data with minimal changes. The data elements collected for HSLs:09 were consistent with items collected for NPSAS:16 (the most recent collection prior to HSLs:09), except that HSLs:09 collected data for six academic years (2011–12 through 2016–17) while NPSAS:16 only covered the 2015–16 academic year.

Some refinements were made to the HSLs:09 student records instrument based on the results of NPSAS:16 data collection and cognitive interviews conducted with participating institutions. These changes were intended to maintain consistency of data elements across NCES postsecondary studies that collect student records data, improve the clarity of item definitions, enhance the usability of the PDP for participating institutions, and facilitate the collection of items across multiple academic years. For example, based on feedback from participating institutions, the instructions for the veteran status item were revised to further distinguish it from the veterans benefits item collected as part of the Financial Aid section.

The HSLs:09 student records instrument could be completed using any combination of three data collection modes:

1. Web mode, in which institution staff used drop-down boxes and text-entry fields to hand-key data directly into the student records instrument, one student at a time within the PDP.
2. Excel mode, in which institutions downloaded a preformatted Excel spreadsheet template from the PDP, keyed or copied student data into a spreadsheet template offline, and then uploaded the completed template to the PDP website.
3. CSV mode, in which institutions downloaded customized file specifications from the PDP website, prepared data files offline according to the file specifications, and then uploaded completed files to the PDP website.

Institutions could choose any of these modes, or use a combination of them, to provide student records data.

Prior to the start of HSLS:09 student records collection, changes were made to the PDP website to improve its usability and the quality of the data collected. For example, the Excel mode template was modified to accept both the pre-formatted response options or a corresponding numeric code, which was provided to users in a codebook document. This change allowed users to more easily copy student records data from other sources and paste them into the Excel template without the need for recoding.

### **3.4.2 Student Records Quality-Control Procedures**

Once institutions submitted their final student records data via the PDP, the data were reviewed for quality and completeness. First, automated programs were used to assess the quality of the data and detect missingness of critical data elements or whether sample members' personally identifying information had been changed by the institution. The automated programs produced a data-quality report, which listed the results of each of the programmatic checks. For each institution, project staff reviewed the data-quality reports, identified the source of any data-quality problems, and if needed, reviewed the student records data submitted by the institution to determine if errors could be resolved by project staff. When data problems could not be resolved, project staff documented the specific data issues in the ICS and sent the school's information to ICs to follow up via telephone. In some cases, the institution indicated that the data should be used "as is" and in others, the institution agreed to provide updated data. The most common reason that institutions were contacted about data-quality problems was because entire data sections or individual critical data elements were missing.

## **3.5 Student Records Data Collection Results**

Of the 3,271 eligible institutions with sampled students, 1,991 institutions (61 percent) provided student records data. About half (51 percent) of the institutions opted for Web mode, 39 percent used Excel mode, and 10 percent uploaded CSV files. Table 4 shows student records collection results by institution type.

**Table 4. Number and percent of participating institutions, by student records collection methods, by institution type: 2018**

Institution type	Total eligible institutions	Institution-level participation <sup>1</sup>		Web mode		Excel mode		CSV mode	
		Number	Percent	Number	Percent	Number	Percent	Number	Percent
<b>Total</b>	<b>3,271</b>	<b>1,991</b>	<b>60.9</b>	<b>1,023</b>	<b>51.4</b>	<b>767</b>	<b>38.5</b>	<b>201</b>	<b>10.1</b>
Institution type									
Public									
Less-than-2-year	55	22	40.0	16	72.7	5	22.7	1	4.5
2-year	787	525	66.7	206	39.2	249	47.4	70	13.3
4-year, non-doctorate-granting	304	193	63.5	68	35.2	94	48.7	31	16.1
4-year, doctorate-granting	334	232	69.5	59	25.4	120	51.7	53	22.8
Private nonprofit									
Less-than-4-year	63	25	39.7	21	84.0	3	12.0	1	4.0
4-year, non-doctorate-granting	514	391	76.1	256	65.5	125	32.0	10	2.6
4-year, doctorate-granting	369	272	73.7	140	51.5	109	40.1	23	8.5
Private for-profit									
Less-than-2-year	311	149	47.9	115	77.2	29	19.5	5	3.4
2-year	258	91	35.3	73	80.2	15	16.5	3	3.3
4-year	216	91	42.1	69	75.8	18	19.8	4	4.4

<sup>1</sup> An institution was considered a participant if it provided information for at least one student. A small number (60) of the participating institutions are not represented in the institution type rows due to unknown institution sector.

NOTE: Detail may not sum to totals because of rounding.

SOURCE: U.S. Department of Education, National Center for Education Statistics, High School Longitudinal Study of 2009 Postsecondary Education Transcript Study and Student Financial Aid Records Collection.

Within the 1,991 institutions that provided student records, 19,254 requests were included. The fielded sample was 17,201 sample members. Of those, 411 were deemed ineligible for student records, leaving 16,790 sample members who were eligible or had unknown eligibility. Aggregated to the student level, a sufficient amount of student records data was received from institutions for 8,688 sample members (52 percent); section 5.1 provides further details on what constituted an adequate quantity of data. Table 5 shows the student records collection results at the student level.

**Table 5. Student-level student records collection results: 2018**

Student sample	Number	Percent
<b>Total fielded sample</b>	<b>17,201</b>	<b>†</b>
Fielded, eligible or unknown eligibility	16,790	97.6
Eligible, student records respondents <sup>1, 2</sup>	8,688	51.7
Eligible, student records nonrespondents <sup>2</sup>	7,644	45.5
Unknown eligibility	458	2.7
Fielded, ineligible	411	2.4

† Not applicable.

<sup>1</sup> A sample member was considered a student records respondent if information on state aid and institution aid awards was provided by the institution that was identified as the student's first primary institution. See section 5.1 for more details.

<sup>2</sup> Of the total sample ( $n = 25,206$ ), 898 sample members were known to be eligible for student records but were not fielded because they were not fielded in the second follow-up.

NOTE: Detail may not sum to totals because of rounding.

SOURCE: U.S. Department of Education, National Center for Education Statistics, High School Longitudinal Study of 2009 Postsecondary Education Transcript Study and Student Financial Aid Records Collection.

## Chapter 4. Data Processing and Editing

The activities taking place between collection of the HSLS:09 PETS and SR data and release of the final source data file products may be coarsely divided into two tasks: processing and editing. During data processing, the data were extracted from SQL tables maintained by the data collection team and subjected to rigorous quality checks to verify data were correctly entered and internally consistent. Following data processing, the data files were edited for clarity and ease of use. This chapter details the data-processing and editing procedures.

### 4.1 Data Processing

The subsections below describe the quality checks performed on the HSLS:09 postsecondary transcripts and student records data.

#### 4.1.1 Transcript Data Reassignment and Consolidation

After data were keyed and coded, as described in section 3.2.2, data were extracted and stored in SAS data files. Transcript data were reviewed daily for inconsistencies and potential keying errors. If keying errors were discovered, the data were corrected within the KCS (see section 3.2.2). In addition to identifying irregularities, project staff performed three main activities to support consistency across transcripts: upcoding, bulk credit review, and calculated variable construction.

**Upcoding.** Several fields within the KCS allowed for “other, specify” text. These fields allowed the system to capture transcript data that did not conform to the specified inputs. When data were entered in these fields, project staff reviewed each character string to determine if it could be upcoded to an existing category. Values that could not be coded into existing categories were left as “other.”

**Bulk credit review.** *Bulk credit* is defined as the total sum of credits a student received for multiple courses or tests. Project staff reviewed instances of bulk credit in an effort to associate portions of the total credit with individual courses or tests. Bulk credit could be decomposed if, for example, a transcript from another institution lists the credits per course (or test) explicitly. For example, suppose a student’s transcript from School B indicates that six credits were accepted from School A. Suppose also that the student’s transcript from School A indicates that three credits were earned for MA 425 and 3 credits earned for COMP 116. In this

case, the six credits accepted at School B would be associated with the two courses listed on School A's transcript. If the credit total could not be decomposed, the credit amount was left in the data as a single bulk sum.

**Calculated variable construction.** The data collected in the KCS may be at such a detailed level that they are not easily analyzed. Calculated variables, which are composite variables at a level other than the student-level (e.g., student-institution-year-level variables), aggregate or combine source data so that the information is more easily accessible. These variables are then stored on the source data files and documented. The calculated variables are created using the same process as student-level composite variables. For details on composite variables, see section 6.5.

### 4.1.2 *Processing Student Records Data*

As discussed in section 3.4.2, upon an institution's submission of data, an initial stage of review verified adherence to the specifications provided. The data collection quality-control team corrected any errors identified. Once an institution's data passed the initial review stage, they were extracted and placed into a SAS dataset for further processing. Project staff performed five main activities to ensure consistency across student records: sanitization, value recoding, financial aid program review, CIP code review, and composite variable construction.

**Sanitization.** Verbatim character strings, such as financial aid strings and major strings, must be sanitized to ensure the integrity of the data and confidentiality of the respondents. Project staff censored character strings provided by institutions by redacting personally identifiable information that could be used to identify respondents.

**Value recoding.** Project staff reviewed the data to ensure that each variable contained valid and consistent values and recoded invalid entries as needed to ensure that data were not lost when submitted in an invalid format. For example, an institution may have indicated a student's enrollment status to be "Full" instead of the requested value, "Full-Time." In such a case, project staff recoded the enrollment status to "Full-Time." This process was executed programmatically such that the status "Full" was only converted once. Consistency with prior student records collections, such as NPSAS:16, was ensured by using a collection of common recodes created during earlier collections.

**Financial aid program review.** Financial aid programs were reviewed to ensure consistent and accurate categorization. For example, if staff identified an aid program by the program name as a state merit grant, but the award was inadvertently

categorized by the institution as an institutional merit grant, the source for the award was changed from “institution” to “state.” This process was also executed programmatically such that aid programs were reviewed once per institution. Consistency with prior student records collections, such as NPSAS:16, was ensured by using a bank of common aid programs created during earlier collections.

**CIP code review.** All major fields of study that contained invalid or blank CIP codes were systematically reviewed using a coding application. Project staff used this application to review the major text field and provide a valid CIP code when possible. This process supported the encoding of valid, consistent, and accurate majors.

**Composite variable construction.** Data collected in the student records instrument may be at such a fine level of detail that they are not easily analyzed. To generate variables that are more easily analyzable, project staff constructed a set of composite variables which aggregate or combine source data so that the information is more accessible. These variables were created using the same process as student-level composite variables. See section 6.5 for further details.

## 4.2 Data Editing, Documentation, and Review

Over the course of data collection, project staff worked to prepare the source data for release on the restricted-use files. This work involved conducting checks of all information collected from institutions to verify the precision and accuracy of data and construction of documentation to aid researchers.

**Data editing.** Routine data inspection is one of the most important quality-control activities performed in data editing. These activities help verify that proper relationships are maintained between variables and that all edits are applied appropriately and consistently. Staff conducted the following steps in editing the data files for release.

- Staff performed logical recoding of the data when the value of missing items could be determined from answers to previous questions. For example, if the institution provided a date of high school completion, but a response to whether the student completed high school was missing, the value for whether the student completed high school was set to “yes,” and the edit was subsequently documented in the codebook.
- Staff assigned labels to the expected values of categorical variables, which aided in revealing any unexpected values. Unexpected values were labeled

when appropriate or set to a reserve code (see section 6.4.1 for further details), in which case the edit was subsequently documented in the codebook.

- Staff examined the minimum, maximum, mean, and median values of continuous variables to assess reasonableness of responses. Staff investigated and corrected or documented anomalous distributions and values. If the value in question was unacceptable, the value was replaced with a reserve code (i.e., -6; see table 21 in section 6.4) and the edit was documented in the codebook.
- Staff examined all missing data to assign specific values indicating the cause of the missing data (see table 21). For example, staff defined gate-nest question<sup>9</sup> relationships and examined data for adherence to logic established in the instrument design, assigning a value of -7 to indicate a legitimate skip.
- Staff generated cross-tabulations of similar and related items to verify that the proper relationships between variables and reserve codes held.

**Documentation.** Accurate and clear data documentation is as important as data editing. Throughout the data collection and data editing process, project staff managed and updated metadata associated with the data files. These metadata consist of variable-by-variable documentation, including variable names, variable labels, value codes, value labels, variable distributions, variable descriptions, conditions under which the variable applies to a particular unit (school or student), and any important notes regarding the creation or use of the variable. All variable documentation was maintained in the Metadata System (MDS), a centralized, browser-based repository. This central system allowed staff to update, refine, and output the latest information while maintaining version control. The MDS was also used to track the status of data revisions throughout the editing process. The team also used this system to generate automated quality-control and progress reports daily.

During data collection, project staff regularly compiled interim datasets for review by NCES project officers and third-party reviewers; upon completion of data collection, a release-ready version of the student records and transcript data files and associated documentation were delivered to NCES for final release.

---

<sup>9</sup> For some questions in the instrument, a question (“gate”) must be first answered before a set of subsequent questions are asked (“nest” or “nested questions”).

## **Chapter 5. Response Rates, Analytic Weights, Variance and Design Effects Estimation, Nonresponse Bias Analysis, Imputation, and Disclosure Avoidance**

The post-data-collection statistical activities conducted to support the analysis of postsecondary transcript and student records data are presented and discussed in this chapter. Section 5.1 describes the criteria for defining respondents to the PETS and SR components of PETS-SR. A discussion of weighted unit response rates from each round of the High School Longitudinal Study of 2009 (HSLs:09) is provided in section 5.2. Section 5.3 includes a succinct description of the weights developed prior to the PETS-SR round and in-depth discussion of the weights developed for PETS-SR. Guidance on the process of selecting weights for particular analyses is provided in section 5.4. The appropriate calculation of standard errors and estimates of the impact of sampling and weight adjustments on the precision of standard errors is discussed in section 5.5. A discussion of bias arising from item nonresponse and unit nonresponse is given in section 5.6 and the methods and results of imputation procedures are presented in section 5.7. Section 5.8 discusses the application of disclosure limitation techniques and explains the resulting differences between public-use and restricted-use data files.

### **5.1 Criteria for Defining Respondents**

For the postsecondary transcripts collection component, a sample member was considered a PETS respondent if at least one transcript was received from an institution attended by the sample member and the transcript contained a course, term, or degree program.

One goal of the student records collection was to report a complete picture of students' financial aid awards throughout their postsecondary enrollments. However, because of the nature of the student records instrument, in which institutions had the ability to answer or skip any item, completeness of the financial aid data varies

significantly across years, institutions, and even across students within institutions. Because the focus of the composite variables was the first primary institution<sup>10</sup> and associated aid packages, to be a student records respondent, information on state aid and institution aid awards must have been provided by the institution that was identified as the student's first primary institution. Federal aid is also a component of the student's total award package, but because these awards are available from the NSLDS, they were not required to have been submitted by the institution. For more information on a student's first primary institution record, see section 6.5.

Of the 19,254 records provided by the contacted institutions, 12,854 were identified to be the student's first primary institution record. Of these first primary institution records, 8,688 had information on state and institution aid awards. Specifically, the first primary institution either provided the amount awarded to the student for at least one state aid program or indicated that the student did not receive any state aid in the first year enrolled, and they either provided the amount awarded to the student for at least one institution aid program or indicated that the student did not receive any institution aid in the first year enrolled.

Although student records composite variables were created for the 8,688 student records respondents only, source student records data are available on the restricted-use file for 205 additional students who are considered respondents to the transcript collection. That is, any student records data that were submitted by the institution for these students are available.

## 5.2 Unit Response Rates

Information on the participation of HSLS:09 sample members is of interest to understand the data collection effort and data quality. Response rates estimate the proportion of the target population represented by sample respondents. The HSLS:09 target population in all rounds prior to PETS-SR is all students in the 9th grade during the fall 2009 term who attended either regular public or private schools in the 50 states and the District of Columbia that provided instruction in both 9th and 11th grades. The PETS target population is a subset of that, as it only consists of

---

<sup>10</sup> The first primary institution is generally the institution in which a student first enrolled at the postsecondary level, according to enrollment data in both transcripts and student records. For students who enrolled at one institution during the summer immediately after high school and enrolled at another institution during the fall, their first primary institution is the institution with the fall enrollment. Academic years are defined as running from July 1 to June 30. The first academic year is generally the earliest academic year in which a student was enrolled at his or her first primary institution. For students who first enrolled in a postsecondary institution in the last 2 months of an academic year and then enrolled the following fall, their first academic year is the academic year of the fall enrollment.

those students in the 9th grade in the fall 2009 term who ever attended an IPEDS-participating postsecondary institution as of June 30, 2017, either before or after completion of high school or a high school equivalency. As discussed in section 2.5, the SR target population is a subset of the PETS population as it consists of students who attended an IPEDS-participating postsecondary institution as of June 30, 2017 after the completion of high school or a high school equivalency. Students who only ever attended a postsecondary institution during high school are ineligible for the SR component. Within this chapter, any reference to a student having attended a postsecondary institution indicates that the student attended an IPEDS-participating postsecondary institution as of June 20, 2017.

In prior rounds of the HSLS:09 study, weighted unit response rates were computed using the base weights. Response rates for the PETS and SR components of PETS-SR were created using the base weights adjusted for sample members whose postsecondary enrollment status was unknown. Details regarding the construction of the weighting adjustments used to account for unknown eligibility are provided in section 5.3. In all rounds, ineligible<sup>11</sup> and sample members known to be deceased at the time of data collection are excluded from the response rate calculation.

As previously mentioned, response rates are used to gauge the degree to which participating schools and participating students represent their respective populations. When response rates are higher, the collected data may produce less biased population estimates, because the larger responding sample may better represent the target population of interest. The weighted unit response rates reported in this DFD report are calculated using the response rate formula provided in NCES Statistical Standard 1-3-2 (Seastrom 2014).

Calculation of a weighted response rate requires identifying the population of interest (school or student) and specifying a participation definition. In studies such as HSLS:09 that are longitudinal in nature and utilize multiple survey components in one or more study round, there are a multitude of participation definitions that may be created. For example, a student participant may be defined as a student who responded to the HSLS:09 second follow-up survey or, alternatively, a student for whom a postsecondary transcript was collected. Response in multiple rounds may be the criteria for participation, where students are considered participants if they responded in the base year and 2013 Update, for example. Several weighted unit

---

<sup>11</sup> Ineligible students are students who never enrolled in a postsecondary course during or after high school completion or receiving a high school equivalency, in the case of PETS, or who never enrolled in a postsecondary course after high school completion or after receiving a high school equivalency, in the case of SR.

response rates, using different definitions of participation and covering all HSLS:09 study rounds, are provided in this section.

Although higher response rates can indicate more accurate survey results, it is also important to examine whether there is the potential for nonresponse bias to exist in the data. NCES standards require unit nonresponse bias analyses to be conducted when weighted unit response rates fall below 85 percent. The base weights account for differential selection probabilities. Analysis weights are constructed by adjusting the base weights for unknown eligibility, if applicable, and nonresponse to mitigate bias induced by those who did not respond to the study. The weights are further calibrated to known population totals to construct analysis weights which enable population estimates to be calculated from sample data.

For some of the survey components in each of the HSLS:09 study rounds, weighted unit response rates computed using the base weights are provided in table 6 as an overview; for a complete listing, see table 10 in section 5.4. Note that schools and students are the sampling units, not parents or teachers; accordingly, response rates are interpreted with respect to schools and students. Note also that there is no specific weight constructed for students who are eligible for, and have, postsecondary transcripts and student records. Weighted response rates incorporating base-year teacher data and rates incorporating multiple sets of data across more than one study round are provided in section 5.4.

**Table 6. HSLs:09 unit response rates**

Unit	Participation definition	Eligible	Participated	Weighted percent
<b>Base Year</b>				
School	School agreed to participate	1,889	944	55.5 <sup>1</sup>
Student	Student questionnaire completed	25,206	21,444	85.7 <sup>2</sup>
	Student assessment completed	25,206	20,781	83.0 <sup>2</sup>
<b>First Follow-up</b>				
Student	Student questionnaire completed <sup>3</sup>	25,184	20,594	82.0 <sup>2</sup>
	Student assessment completed <sup>3</sup>	25,184	18,507	73.0 <sup>2</sup>
	Parent questionnaire completed <sup>5</sup>	11,952	8,621	72.5 <sup>4</sup>
<b>2013 Update and High School Transcript components</b>				
Student	Student questionnaire completed	25,168	18,558	73.1 <sup>2</sup>
	High school transcripts collected	25,167	21,928	87.7 <sup>2</sup>
	Student questionnaire completed and high school transcripts collected	25,167	17,656	70.2 <sup>2</sup>
<b>Second Follow-up</b>				
Student	Student questionnaire completed	25,123	17,335	67.9 <sup>2</sup>
<b>Postsecondary Education Transcript Study and Student Financial Aid Records Collection</b>				
Student	Postsecondary transcript collected	17,338	13,160	71.2 <sup>6</sup>
	Postsecondary student records collected	17,230	8,688	48.7 <sup>7</sup>

<sup>1</sup> Weighted percentage is calculated using the school base weight.

<sup>2</sup> Weighted percentages are calculated using the student base weight.

<sup>3</sup> A total of 22 students from the base year were ineligible for the first follow-up.

<sup>4</sup> Weighted percentage is calculated using the student base-weight adjustment for parent subsampling.

<sup>5</sup> A subsample of 11,952 eligible parents were asked to participate in the HSLs:09 first follow-up data collection.

<sup>6</sup> Weighted percentage is calculated using the student base weight adjusted for unknown eligibility with respect to the Postsecondary Education Transcript Study.

<sup>7</sup> Weighted percentage is calculated using the student base weight adjusted for unknown eligibility with respect to the Student Financial Aid Records Collection.

NOTE: There is no student base weight adjusted for eligibility to both student records collection and postsecondary transcript collection.

SOURCE: U.S. Department of Education, National Center for Education Statistics, High School Longitudinal Study of 2009 (HSLs:09) Postsecondary Education Transcript Study and Student Financial Aid Records Collection, Public-use Data File.

As shown in table 6, weighted response rates for the student questionnaire, which is the only component included in all four data collections, ranged from 85.7 percent in the base year to 67.9 percent in the second follow-up. The weighted response rates of the two postsecondary components vary with 71.2 percent in the postsecondary transcript collection and 48.7 percent achieved for student records collection.

## 5.3 Overview of Weighting

This section describes the purpose of weighting and when weights ought to be incorporated into analyses. An overview is given of the two methods of variance estimation supported with HSLS:09 weights, namely balanced repeated replication (BRR) and Taylor series linearization. The weights which were created in all prior rounds of HSLS:09 are listed, with references to prior rounds' documentation where the readers can obtain further details. The four weights constructed for the PETS-SR study round are discussed in detail.

### 5.3.1 Analysis Weights

The use of weights is essential to produce estimates that are representative of the HSLS:09 target population of students for each study round and component. An analysis weight should be used to produce survey estimates. When testing hypotheses (e.g., conducting  $t$  tests and regression analyses) using weighted data from a study such as HSLS:09 that has a complex design, analysts also should use methods to properly estimate variances. Variables have been created for HSLS:09 to support two methods of variance estimation that account for the HSLS:09 complex sample design: (1) a BRR variance estimation method using the BRR weights and the associated analysis weight and (2) a linearization variance estimation method through a Taylor series approximation using analysis weights and variables that represent school sampling strata and primary sampling units. For more details on standard error estimation, see section 5.5.

#### 5.3.1.1 Weighting in the base year, first follow-up, 2013 Update, and second follow-up

Five sets of weights were constructed for the HSLS:09 base year. The steps implemented to create the five weights are detailed in the *HSLS:09 Base-Year Data File Documentation* (Ingels et al. 2011). The five weights are designed for the following analyses:

- W1SCHOOL: school-level analyses of information collected in the administrator and counselor questionnaires, as well as school-level data from other sources, such as CCD and PSS;
- W1STUDENT: student-level analyses of student survey responses and mathematics assessment scores;
- W1SCITCH: student-level analyses of science teacher questionnaire data;

- W1MATHCH: student-level analyses of math teacher questionnaire data; and
- W1PARENT: student-level analyses of parent questionnaire data.

Four sets of weights were computed for the HSLS:09 first follow-up. The steps utilized to create these weights are discussed in detail in the *HSLS:09 Base Year to First Follow-Up Data File Documentation* (Ingels et al. 2013). The four weights are designed for the following student-level analyses:

- W2STUDENT: analyses specific to the first follow-up student survey;
- W2W1STU: analyses examining both base-year and first follow-up student survey data;
- W2PARENT: analyses specific to data from the first follow-up parent questionnaire; and
- W2W1PAR: analyses examining both base-year and first follow-up parent survey data.

Eleven sets of weights—five High School Transcript weights and six nontranscript weights—were computed for the HSLS:09 2013 Update and High School Transcript study. The steps used to construct four of the nontranscript weights are detailed in the *HSLS:09 2013 Update and High School Transcript Data File Documentation* (Ingels et al. 2015). The following four nontranscript weights are designed for the following student-level analyses:

- W3STUDENT: analyses specific to the 2013 Update;
- W3W1STU: analyses examining both base-year and 2013 Update data;
- W3W2STU: analyses examining both the first follow-up and the 2013 Update; and
- W3W1W2STU: analyses examining base-year, first follow-up, and the 2013 Update data.

Two additional 2013 Update nontranscript weights were constructed during a subsequent round. These weights, and the steps used to create them, are detailed in the *HSLS:09 Base-Year to Second Follow-Up Data File Documentation* (Duprey et al. 2018). These weights are designed for the following student-level analyses:

- W3W1MATHTCH: analyses examining base-year student and math teacher data, and 2013 Update data; and
- W3W1SCITCH: analyses examining base-year student and science teacher data, and 2013 Update data.

The steps used to construct the five High School Transcript weights are detailed in the *HSLs:09 2013 Update and High School Transcript Data File Documentation* (Ingels et al. 2015). The five High School Transcript weights are designed for the following student-level analyses:

- W3HSTRANS: analyses specific to High School Transcript data only;
- W3STUDENTTR: analyses examining 2013 Update data combined with High School Transcript data;
- W3W1STUTR: analyses examining base-year, 2013 Update, and High School Transcript data;
- W3W2STUTR: analyses examining first follow-up, 2013 Update, and High School Transcript data; and
- W3W1W2STUTR: analyses examining base-year, first follow-up, 2013 Update, and High School Transcript data.

Five sets of weights were computed for the HSLs:09 second follow-up. The steps utilized to create these weights are discussed in detail in the *HSLs:09 Base-Year to Second Follow-Up Data File Documentation* (Duprey et al. 2018). The five weights are designed for the following student-level analyses:

- W4STUDENT: analyses specific to the second follow-up;
- W4W1STU: analyses examining both base-year and second follow-up data;
- W4W1W2W3STU: analyses examining base-year, first follow-up, 2013 Update, and second follow-up data;
- W4W1STUP1: analyses examining base-year student and parent data and second follow-up data; and
- W4W1STUP1P2: analyses examining base-year student and parent data, first follow-up parent data, and second follow-up data.

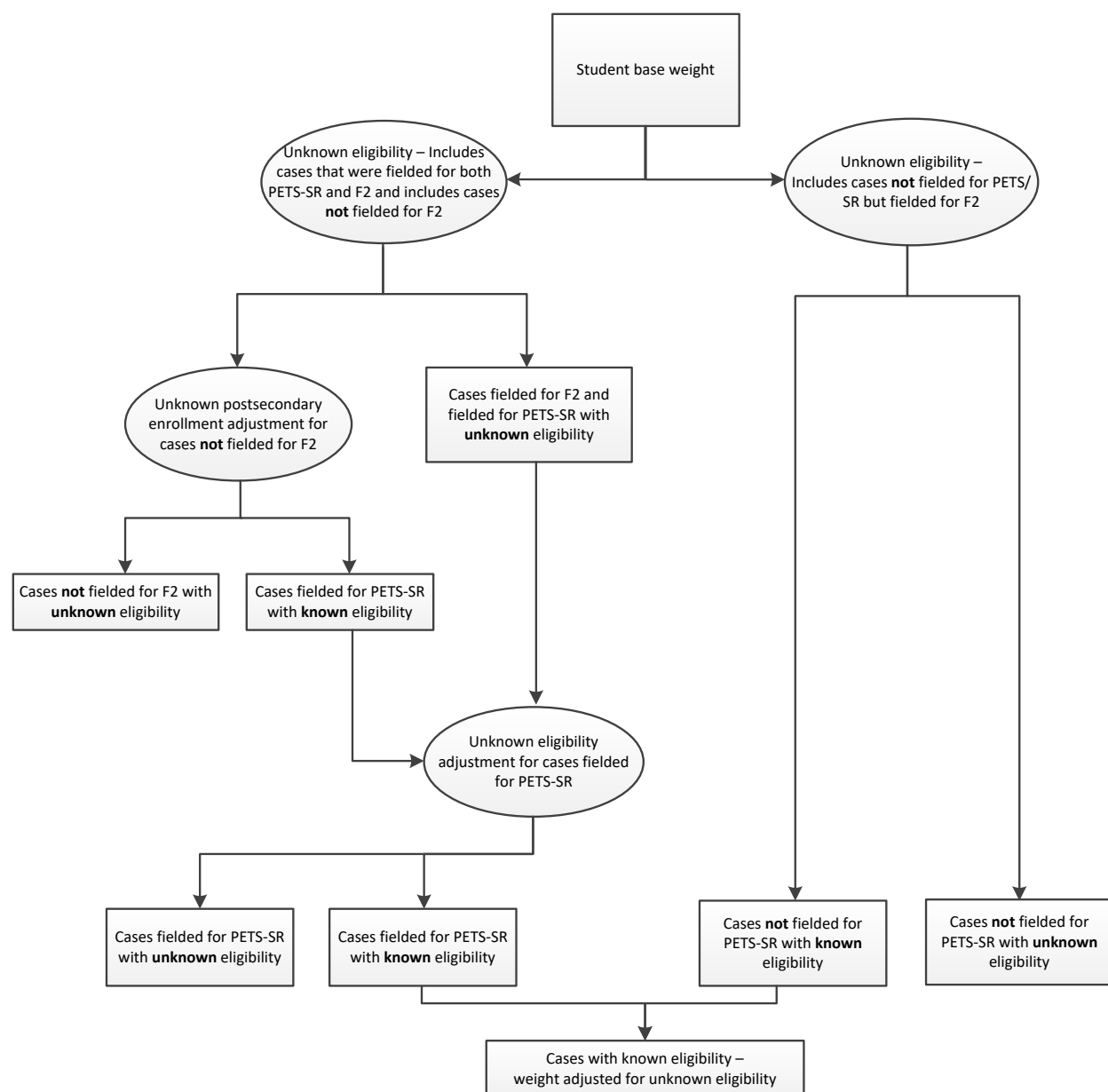
### **5.3.1.2 Weighting in the Postsecondary Education Transcript Study and Student Financial Aid Records Collection**

Four analysis weights were constructed in the PETS-SR study round. The weight names and supported analyses are

- W5PSTRANS: analyses specific to postsecondary transcript data;
- W5W1W2W3W4PSTRANS: analyses examining student survey data from the base-year, first follow-up, 2013 Update, and second follow-up and postsecondary transcript data;
- W5PSRECORDS: analyses specific to postsecondary student records data; and
- W5W1W2W3W4PSRECORDS: analyses examining student survey data from the base-year, first follow-up, 2013 Update, second follow-up and postsecondary student records data.

The weighting adjustments for all four analysis weights followed the same sequence: beginning with adjustments for unknown eligibility, then nonresponse, and finally a calibration to known control totals. When a sample member is neither known to be eligible nor known to be ineligible, they are said to have unknown eligibility. In these instances, unknown eligibility weighting adjustments may be appropriate, such as in HSLS:09 PETS-SR. In an unknown eligibility adjustment, the weight of the sample members with unknown eligibility is distributed to the cases with known eligibility. The resulting weight for cases with known eligibility is said to be adjusted for unknown eligibility.

The adjustments for unknown eligibility followed the same form for all four weights. However, the adjustments themselves are not the same because, as discussed in section 5.2, the definitions for eligibility differ between PETS and SR. Thus, the set of eligible cases and the set of cases with unknown eligibility differ with respect to the PETS and SR weights, requiring one unknown eligibility adjustment for the two PETS weights, and another for the two SR weights. The form of the unknown eligibility adjustments is displayed in figure 6.

**Figure 6. PETS-SR unknown eligibility adjustment construction: 2018**

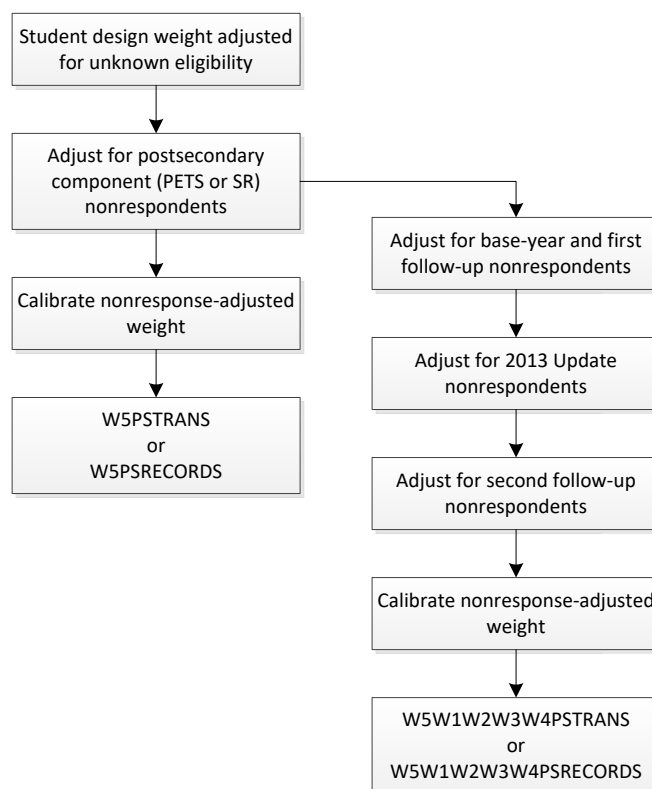
NOTE: F2 = Second follow-up of HSLS:09.

SOURCE: U.S. Department of Education, National Center for Education Statistics, High School Longitudinal Study of 2009 Postsecondary Education Transcript Study and Student Financial Aid Records Collection.

The sample members were split into two different paths for the unknown eligibility adjustment. The first path included two sets of students: (1) students fielded for PETS-SR as described in section 2.5 and (2) students not fielded for the second follow-up as described in section 2.4. Note that all students not fielded for the second follow-up were also not fielded for PETS-SR. The second path included all students not fielded for PETS-SR who were fielded for the second follow-up, as described in sections 2.4 and 2.5. The students were separated based on their fielded status for PETS-SR because differences in key characteristics were observed between cases fielded and not fielded for PETS-SR. The separation through the course of the unknown eligibility adjustments allowed for the distribution of those characteristics to be preserved within each PETS-SR fielded status in the weight adjusted for unknown eligibility.

The cases not fielded for the second follow-up were combined in the path with students fielded for PETS-SR because the cases not fielded for the second follow-up more closely resembled the cases fielded for PETS-SR rather than the other cases not fielded for PETS-SR on the key characteristics evaluated. The cases not fielded for the second follow-up all had unknown eligibility with respect to PETS and SR, and thus all cases not fielded for the second follow-up distribute their weight to cases fielded for PETS-SR with known eligibility. The rationale for a separate unknown eligibility adjustment for cases not fielded for the second follow-up is twofold. First, cases not fielded for the second follow-up were not given the chance to have their eligibility for PETS or SR determined, thus they all have unknown eligibility for the same reason. Second, the majority of cases not fielded for the second follow-up were not fielded in the second follow-up because they never responded in an HSLS:09 study round, see sections 2.3 and 2.4, and thus very limited information was available to inform a model for unknown eligibility involving the cases not fielded for the second follow-up.

The cases fielded for PETS-SR with known eligibility had their base weight completely adjusted for unknown eligibility at the conclusion of the two adjustments in the first path. The cases not fielded for PETS-SR with known eligibility had their base weight adjusted for unknown eligibility once the single adjustment in the second path is completed. All cases with known eligibility were then combined, forming the set of cases with known eligibility and a weight adjusted for unknown eligibility. This weight adjusted for unknown eligibility was the weight input into the nonresponse adjustments displayed in figure 7 and discussed below. Note that only eligible cases were included in the nonresponse adjustments.

**Figure 7. PETS-SR nonresponse and calibration weighting adjustment construction: 2018**

NOTE: PETS = Postsecondary Education Transcript Study; SR = Student Financial Aid Records.

SOURCE: U.S. Department of Education, National Center for Education Statistics, High School Longitudinal Study of 2009 Postsecondary Education Transcript Study and Student Financial Aid Records Collection.

Similar to the unknown eligibility adjustments, the nonresponse adjustments for the two PETS weights and the two SR weights followed the same patterns, but the adjustments themselves were constructed independently for the PETS and for the SR weights. Specifically, the first nonresponse adjustment for each weight was an adjustment for nonresponse to PETS or SR, whichever component was applicable. The adjustment for nonresponse to PETS was the sole nonresponse adjustment for W5PSTRANS and the adjustment for nonresponse to SR was the sole nonresponse adjustment for W5PSRECORDS. To create the weights that account for student nonresponse in all four prior rounds, the weight adjusted PETS or SR nonresponse was further adjusted in three phases. First, it was adjusted for student nonresponse to the base year or first follow-up, then adjusted for nonresponse in the 2013 Update, and finally adjusted for nonresponse in the second follow-up. The weight adjustments for nonresponse to PETS and further adjustments for student nonresponse in the four rounds were the nonresponse adjustments for W5W1W2W3W4PSTRANS. The weight adjustments for nonresponse to SR and

further adjustments for student nonresponse in the four rounds were the nonresponse adjustments for W5W1W2W3W4PSRECORDS.

At the conclusion of the nonresponse adjustment(s) the nonresponse-adjusted weight was calibrated to the same control totals defined in the base year and used in all prior rounds of HSLs:09. However, these control totals are representative of all students in the 9th grade during the fall 2009 term who attended either regular public or private schools in the 50 states and the District of Columbia that provided instruction in both 9th and 11th grades. As discussed in section 5.2, the population that the PETS weights represent is a subset of that population, and the population that the SR weights represent is a further subset of the PETS population. Thus, deceased and ineligible<sup>12</sup> students were included in the calibration and their weights were subsequently set to zero. The resulting two calibrated PETS weights, W5PSTRANS and W5W1W2W3W4PSTRANS, are representative of the portion of the HSLs:09 base-year population that is not deceased and ever enrolled in a postsecondary institution. The two calibrated SR weights, W5PSRECORDS and W5W1W2W3W4PSRECORDS, are representative of the portion of the HSLs:09 base-year population that is not deceased and enrolled in a postsecondary institution after the completion of high school or a high school equivalency.

For all four PETS-SR weights, unknown eligibility and unit nonresponse adjustments incorporated student-level and school-level characteristics where possible using the WTADJUST procedure in SUDAAN. The calibrations for each weight also used the WTADJUST procedure in SUDAAN.

For more detail on the construction of all four PETS-SR weights, please refer to appendix E. Additional information on using the analysis weights to estimate standard errors is provided in section 5.5.1.

### 5.3.2 BRR Weights

A set of 200 BRR weights was created for each of the four PETS-SR analysis weights. These sets of BRR weights included (1) postsecondary transcript student weights (W5PSTRANS001–200); (2) base-year to first follow-up to 2013 Update to second follow-up with a postsecondary transcript student weights (W5W1W2W3W4PSTRANS001–200); (3) postsecondary student records student weights (W5PSRECORDS001–200); and (4) base-year to first follow-up to 2013 Update to second follow-up with postsecondary student records student weights (W5W1W2W3W4PSRECORDS001–200). Procedures for constructing the weights

---

<sup>12</sup> All students ineligible for PETS are ineligible for SR. Additional students who are eligible for PETS are ineligible for SR.

mirrored those used to construct the corresponding analysis weight. Namely, the BRR weights were constructed by subjecting the base-year BRR base weights, defined for each of 200 replicates, to unknown eligibility, nonresponse, and calibration adjustments following a process like that used to develop the analysis weights. Additional information on using the BRR weights to estimate standard errors may be found in section 5.5.1.

### 5.3.3 Weight Characteristics

The characteristics of the four PETS-SR analysis weights are presented in table 7. For each weight, the number of respondents, the mean weight, the standard deviation, the minimum and maximum, and weight sums are provided.

**Table 7. Descriptive characteristics of PETS-SR survey weights: 2018**

Weight	Number of respondents	Mean	Standard deviation	Minimum	Maximum	Sum <sup>1</sup>
W5PSTRANS	13,160	248.6	327.18	2.0	6,633.0	3,271,346
W5W1W2W3W4PSTRANS	8,932	363.5	495.20	4.6	7,174.7	3,247,165
W5PSRECORDS	8,688	373.9	488.54	4.9	6,720.9	3,248,410
W5W1W2W3W4PSRECORDS	5,936	544.1	749.94	6.9	9,932.7	3,229,723

<sup>1</sup> The student counts in table 10 of chapter 3 in the *HSLS:09 Base-Year Data File Documentation* (Ingels et al. 2011) were used as the control totals. Weight sums differ from the population counts for three reasons: (1) suppression of data from the public-use file for the students who were excluded from the base-year or first follow-up student survey because it was not offered in a format that allowed their meaningful participation (students referred to as “questionnaire incapable” in response status variables) in the base year or first follow-up; (2) deceased students were included in the calibration and subsequently had their weights set to zero; and (3) students who were ineligible with respect to PETS or SR collection were included in the calibration and subsequently had their weights set to zero.

SOURCE: U.S. Department of Education, National Center for Education Statistics, High School Longitudinal Study of 2009 (HSLS:09) Postsecondary Education Transcript Study and Student Financial Aid Records Collection, Public-use Data File.

The weighted counts and percentages of the X2SEX variable for the restricted- and public-use file by the four PETS-SR analysis weights are presented in table 8.

**Table 8. Weighted counts and percentages of X2SEX for restricted- and public-use files, by PETS-SR survey weight: 2018**

Weight	Restricted-use file				Public-use file			
	Males		Females		Males		Females	
	Count	Percent	Count	Percent	Count	Percent	Count	Percent
W5PSTRANS	1,562,173	47.75	1,709,173	52.25	1,562,173	47.75	1,709,173	52.25
W5W1W2W3W4PSTRANS	1,553,756	47.85	1,693,409	52.15	1,553,756	47.85	1,693,409	52.15
W5PSRECORDS	1,545,820	47.59	1,702,590	52.41	1,545,820	47.59	1,702,590	52.41
W5W1W2W3W4PSRECORDS	1,541,380	47.72	1,688,343	52.28	1,541,380	47.72	1,688,343	52.28

NOTE: PETS = Postsecondary Education Transcript Study; SR = Student Financial Aid Records. Weighted counts and percentages do not incorporate students who were excluded from the base-year or first follow-up student survey because it was not offered in a format that allowed their meaningful participation (students referred to as "questionnaire incapable" in response status variables) in the base year or first follow-up. Counts and percentages are weighted by the weight in the row of interest within the file of the given column.

SOURCE: U.S. Department of Education, National Center for Education Statistics, High School Longitudinal Study of 2009 (HSLs:09) Postsecondary Education Transcript Study and Student Financial Aid Records Collection, Restricted-use and Public-use Data Files.

### 5.3.4 Weighting Quality Control

A good weight is one which allows an analyst to adequately adjust for unknown eligibility, unit nonresponse, and adjust for frame coverage through calibration while minimizing overall weight variability. To assess the quality of the analysis weights and their corresponding sets of replicate weights, staff reviewed the following:

- the initial base weights' characteristics, including the (1) distribution of the weights, (2) ratios of maximum weights to minimum weights, (3) unequal weighting effects, (4) ranges of weight adjustment factors, and (5) weight sums;
- the weight adjustment factors used to produce the postsecondary transcript and student records weights; and
- the variability of the weights themselves and the degree to which the sums of the individual weights matched calibration totals.

Some of the specific quality-control checks employed for unknown eligibility, nonresponse, and calibration weight adjustments are described below.

Unknown eligibility and nonresponse weight adjustment quality-control checks included the following:

- Weight sums after the adjustment matched weight sums before the adjustment. This assessment included the overall weight sum and weights sums by the levels of the categorical variables used in the weighting model, such as student race/ethnicity, sex, and base-year school type.
- Overall unequal weighting effect (UWE) after weighting adjustment was not substantially higher than the overall UWE prior to weighting adjustment. As a general rule of thumb, increases in the overall UWE were kept within 10 percent of the overall UWE prior to nonresponse adjustment.
- UWEs before nonresponse adjustment were computed for the main effect variables<sup>13</sup> used in the nonresponse models and compared to the corresponding UWEs after nonresponse adjustment. As a general rule of thumb, increases of 10 percent in the UWEs were considered acceptable.

---

<sup>13</sup> Main effect variables are variables such as sex and race/ethnicity, which are distinguished from interacted variables such as school type within Census region.

Calibration weight adjustments quality-control checks included the following:

- Weight sums after calibration were compared with target control totals to verify equivalence. The target control totals included totals defined for the school type, region, state (if part of the 10 state-representative public-school samples), and metropolitan status.
- UWEs were computed before calibration and compared with corresponding UWEs after calibration by school type, region, augmented state, metropolitan status, sex, and race.
- Design effects for key variables within the associated component for the weight were evaluated to ensure that the design effects are of acceptable size, such as achieving a root design effect (*deft*) below five for as many student domains as possible.

Additional quality-control checks for BRR weight construction included the following:

- Comparing overall UWEs, minimum weights, maximum weights, and average weights across each set of 200 replicates to verify comparability of the replicate weight distributions.
- Design effects for key variables within the associated component for the weight were evaluated and compared to the design effects of the corresponding analysis weight to ensure they are equivalent.

## 5.4 Choosing an Analytic Weight

The choice about which weights to create for HSLS:09 data is driven by the need to maximize the analysis utility for the research community. Analyses may incorporate data obtained from a particular instrument within a round of the study (e.g., postsecondary transcripts in PETS-SR) or combinations of data from multiple instruments across multiple rounds, such as student and parent questionnaire responses in the base year and first follow-up. As discussed in the *HSLS:09 Base-Year to Second Follow-Up Data File Documentation* (Duprey et al. 2018) and repeated here,

weights were derived that incorporate many, but not all, possible combinations of data sources and rounds of data collection.<sup>14</sup>

The PETS-SR data file contains a total of 29 analysis weights: five weights for analysis of the base-year data, four weights to be used in conjunction with the first follow-up data, six weights to be used for analysis involving the 2013 Update, five weights for analyses using High School Transcript data, five weights for analyses of second follow-up data, and four weights for the analysis of student postsecondary transcript and student records data. Guidance for use of all 29 analysis weights is provided in table 9.

The analysis weights presented in table 9 can be used for analysis of data collected in a single study round or data collected across multiple study rounds. The weights designed to be used in analysis of a single round of data are classified as “single-round” weights, and the weights that may be used to analyze data collected from multiple study rounds are classified as “multiround.”

Analyses of base-year data involving only the student assessment data or student questionnaire responses should use W1STUDENT, and base-year analyses that include parent responses from the base year should utilize W1PARENT. Analysis of school administrator or counselor responses, in the context of the HSLS:09 base-year school population, should utilize W1SCHOOL. Similarly, analyses involving only the first follow-up student questionnaire or assessment data should utilize W2STUDENT, and analyses involving first follow-up parent responses should utilize W2PARENT. Analyses involving only 2013 Update data should use W3STUDENT for analyzing questionnaire responses, and analyses of only High School Transcript data should use W3HSTRANS. Analyses involving 2013 Update questionnaire and High School Transcript data only should use W3STUDENT\*TR. Analyses that involve only second follow-up questionnaire responses should use W4STUDENT. Any analysis only involving postsecondary transcript data should use W5PSTRANS and analyses that involves only postsecondary student records data should use W5PSRECORDS.

Some of the analysis weights presented in table 9 are appropriate to use when analyzing data collected across multiple study rounds or from multiple data sources. For example, an analysis seeking to determine base-year predictors of on-time high school graduation should incorporate the analysis weight W3W1STUTR. Similarly, an analysis seeking to determine prior-round predictors of income as of the second

---

<sup>14</sup> The creation of additional HSLS:09 weights was considered. However, to limit potential confusion in the choice of analysis weight if a large number of weights were produced, decisions were made to focus only on the most likely types of analyses given the HSLS:09 data sources.

follow-up should use  $W4W1W2W3STU$  if the set of possible predictors was limited to student questionnaire responses. If postsecondary student records data was added to that analysis of prior-round data, then  $W5W1W2W3W4PSRECORDS$  is the recommended weight.

**Table 9. HSLs:09 analysis weights: 2018**

<b>HSLs:09 round(s)</b>	<b>Universe<sup>1</sup></b>	<b>Estimation</b>	<b>Variable name</b>	<b>Nonresponse-adjusted component(s) in each weight<sup>2</sup></b>
Base year	All study-eligible schools	Single-round	W1SCHOOL	BY School
Base year	All study-eligible students in base year <sup>3</sup>	Single-round	W1STUDENT	BY Student
			W1PARENT	BY Student & BY Parent
			W1SCITCH	BY Student & BY Science teacher
			W1MATHTCH	BY Student & BY Math teacher
First follow-up	9th-grade cohort <sup>3</sup>	Single-round	W2STUDENT	F1 Student
			W2PARENT	F1 Parent
Base year and first follow-up	9th-grade cohort <sup>3</sup>	Multiround	W2W1STU	BY/F1 Student
			W2W1PAR	BY/F1 Student & BY/F1 Parent
2013 Update	9th-grade cohort <sup>4</sup>	Single-round	W3STUDENT	U13 Student
Base year and 2013 Update	9th-grade cohort <sup>3,4</sup>	Multiround	W3W1STU	BY/U13 Student
			W3W1MATHTCH	BY/U13 Student & BY Math teacher <sup>5</sup>
			W3W1SCITCH	BY/U13 Student & BY Science teacher <sup>6</sup>
First follow-up and 2013 Update	9th-grade cohort <sup>3,4</sup>	Multiround	W3W2STU	F1/U13 Student
Base year, first follow-up, and 2013 Update	9th-grade cohort <sup>3,4</sup>	Multiround	W3W1W2STU	BY/F1/U13 Student
High School Transcript	9th-grade cohort <sup>4</sup>	Single-round	W3HSTRANS	High School Transcript

See notes at end of table.

**Table 9. HSLs:09 analysis weights: 2018—continued**

<b>HSLs:09 round(s)</b>	<b>Universe<sup>1</sup></b>	<b>Estimation</b>	<b>Variable name</b>	<b>Nonresponse-adjusted component(s) in each weight<sup>2</sup></b>
High School Transcript and 2013 Update	9th-grade cohort <sup>4</sup>	Single-round	W3STUDENTTR	High School Transcript & U13 Student
High School Transcript, base year, and 2013 Update	9th-grade cohort <sup>3,4</sup>	Multiround	W3W1STUTR	High School Transcript & BY/U13 Student
High School Transcript, base year, first follow-up, and 2013 Update	9th-grade cohort <sup>3,4</sup>	Multiround	W3W1W2STUTR	High School Transcript & BY/F1/U13 Student
High School Transcript, first follow-up, and 2013 Update	9th-grade cohort <sup>3,4</sup>	Multiround	W3W2STUTR	High School Transcript & F1/U13 Student
Second follow-up	9th-grade cohort <sup>4</sup>	Single-round	W4STUDENT	F2 Student
Base year and second follow-up	9th-grade cohort <sup>3,4</sup>	Multiround	W4W1STU	BY/F2 Student
			W4W1STUP1	BY/F2 Student & BY Parent <sup>7</sup>
Base year, first follow-up, and second follow-up	9th-grade cohort <sup>3,4</sup>	Multiround	W4W1STUP1P2	BY/F2 Student & BY/F1 Parent <sup>8</sup>
Base year, first follow-up, 2013 Update, and second follow-up	9th-grade cohort <sup>3,4</sup>	Multiround	W4W1W2W3STU	BY/F1/U13/F2 Student
Postsecondary transcript	9th-grade cohort ever enrolled in a postsecondary institution <sup>4</sup>	Single-round	W5PSTRANS	Postsecondary transcript
Postsecondary transcript, base year, first follow-up, 2013 Update, and second follow-up	9th-grade cohort ever enrolled in a postsecondary institution <sup>3,4</sup>	Multiround	W5W1W2W3W4PSTRANS	Postsecondary transcript & BY/F1/U13/F2 Student

See notes at end of table.

**Table 9. HSLs:09 analysis weights: 2018—continued**

HSLs:09 round(s)	Universe <sup>1</sup>	Estimation	Variable name	Nonresponse-adjusted component(s) in each weight <sup>2</sup>
Postsecondary student records	9th-grade cohort ever enrolled in a postsecondary institution after completion of high school or equivalency <sup>4</sup>	Single-round	W5PSRECORDS	Postsecondary student records
Postsecondary student records, base year, first follow-up, 2013 Update, and second follow-up	9th-grade cohort ever enrolled in a postsecondary institution after completion of high school or equivalency <sup>3, 4</sup>	Multiround	W5W1W2W3W4PSRECORDS	Postsecondary student records & BY/F1/U13/F2 Student

<sup>1</sup> The sum of the associated analysis weights estimates the population count for the universe.

<sup>2</sup> Student-level weights are derived from the school analysis weight and therefore are also adjusted for school nonresponse. Unless otherwise specified, the weights were additionally adjusted for nonresponse within the specified round(s) of data collection.

<sup>3</sup> In the public-use student files, the student weights that require student participation in the base year, first follow-up, or both, generalize to the base-year study-eligible students who would have been capable of completing the corresponding student questionnaire and math assessment in the base year, first follow-up, or both rounds, depending on the rounds used to construct the weight. For illustration, in the student public-use files, W2STUDENT generalizes to the base-year study-eligible students who were alive as of the first follow-up who were able to complete the first follow-up student questionnaire and math assessment. In the case of W2W1STU, this weight generalizes to the base-year study-eligible students who were alive as of the first follow-up who were able to complete the student questionnaire and math assessment in the base-year and the first follow-up.

<sup>4</sup> Excludes those from the cohort who were deceased at the time of the latest data collection accounted for by the weight.

<sup>5</sup> Accounts for student nonresponse in the base year, nonresponse in the 2013 Update, and base-year math teacher nonresponse.

<sup>6</sup> Accounts for student nonresponse in the base year, nonresponse in the 2013 Update, and base-year science teacher nonresponse.

<sup>7</sup> Accounts for student nonresponse in the base year, nonresponse in the second follow-up, and base-year parent nonresponse.

<sup>8</sup> Accounts for student nonresponse in the base year, nonresponse in the second follow-up, base-year parent nonresponse, and first follow-up parent nonresponse.

NOTE: BY = base year; F1 = first follow-up; F2 = second follow-up; U13 = 2013 Update.

SOURCE: U.S. Department of Education, National Center for Education Statistics, High School Longitudinal Study of 2009 (HSLs:09) Postsecondary Education Transcript Study and Student Financial Aid Records Collection.

The number and percentage of completed surveys, high school transcript responses, postsecondary transcript and student records responses, or their combinations for the student sample, and associated recommended weights for the HSLs:09 base-year, first follow-up, 2013 Update, High School Transcript, second follow-up, and PETS-SR study rounds are summarized in table 10. Please note that, although the restricted-use file contains nonzero weights for students who were excluded from the base-year or first follow-up student survey because it was not offered in a format that allowed their meaningful participation (students referred to as “questionnaire incapable” in response status variables) in the base year or first follow-up, the weights for such students are set to zero in the corresponding public-use files.

**Table 10. Number and percentage of completed surveys, high school transcript responses, postsecondary transcript and student records responses, or their combinations for the student sample, and associated recommended weights: PETS-SR**

Study round and high school transcript combinations	Data source(s) and recommended weights	Eligible	Participated	Weighted percent <sup>1</sup>	Unweighted percent
Base year	BY Student questionnaire (W1STUDENT <sup>2</sup> )	25,206	21,444	85.7	85.1
	BY Student assessment (W1STUDENT <sup>3</sup> )	25,206	20,781	83.0	82.4
	BY Student and Parent questionnaires (W1PARENT <sup>2</sup> )	25,206	16,429	65.3	65.2
	BY School administrator (W1STUDENT <sup>3</sup> )	25,206	20,301	81.1	80.5
	BY School counselor (W1STUDENT <sup>3</sup> )	25,206	19,505	77.7	77.4
	BY Teacher questionnaire <sup>4</sup>				
	Math teacher (W1MATHTECH <sup>2</sup> )	23,621	16,035	65.1	67.9
	Science teacher (W1SCITCH <sup>2</sup> )	22,597	14,629	63.6	64.7
First follow-up	F1 Student questionnaire (W2STUDENT <sup>2</sup> )	25,184	20,594	82.0	81.8
	F1 Student assessment (W2STUDENT <sup>3</sup> )	25,184	18,507	73.0	73.5
	F1 Parent questionnaire <sup>5</sup> (W2PARENT <sup>2</sup> )	11,952	8,621	72.5	72.4
Base year and first follow-up	BY/F1 Student questionnaires (W2W1STU <sup>2</sup> )	25,184	18,623	74.3	74.0
	BY/F1 Student assessments (W2W1STU <sup>3</sup> )	25,184	16,356	64.7	65.0
	BY/F1 Student and Parent questionnaires <sup>6</sup> (W2W1PAR <sup>2</sup> )	11,952	6,371	52.9	53.3
2013 Update	U13 Student questionnaire (W3STUDENT <sup>2</sup> )	25,168	18,558	73.1	73.7

See notes at end of table.

**Table 10. Number and percentage of completed surveys, high school transcript responses, postsecondary transcript and student records responses, or their combinations for the student sample, and associated recommended weights: PETS-SR—Continued**

Study round and high school transcript combinations	Data source(s) and recommended weights	Eligible	Participated	Weighted percent <sup>1</sup>	Unweighted percent
Base year and 2013 Update	BY/U13 Student questionnaires (W3W1STU <sup>2</sup> )	25,168	17,117	67.6	68.0
	BY/U13 Student and BY Teacher questionnaires				
	Math teacher <sup>7</sup> (W3W1MATHTCH <sup>2</sup> )	23,587	12,812	51.4	54.3
	Science teacher <sup>8</sup> (W3W1SCITCH <sup>2</sup> )	22,566	11,803	50.7	52.3
First follow-up and 2013 Update	F1/U13 Student questionnaires (W3W2STU <sup>2</sup> )	25,168	17,282	68.0	68.7
Base year, first follow-up, and 2013 Update	BY/F1/U13 Student questionnaires (W3W1W2STU <sup>2</sup> )	25,168	15,857	62.5	63.0
High School Transcript	High School Transcript (W3HSTRANS <sup>2</sup> )	25,167	21,928	87.7	87.1
High School Transcript and 2013 Update	High School Transcript and U13 Student questionnaire (W3STUDENTTR <sup>2</sup> )	25,167	17,656	70.2	69.6
High School Transcript, base year, and 2013 Update	High School Transcript and BY/U13 Student questionnaires (W3W1STUTR <sup>2</sup> )	25,167	16,303	64.7	64.4
High School Transcript, first follow-up, and 2013 Update	High School Transcript and F1/U13 Student questionnaires (W3W2STUTR <sup>2</sup> )	25,167	16,525	65.6	64.9
High School Transcript, base year, first follow-up, and 2013 Update	High School Transcript and BY/F1/U13 Student questionnaires (W3W1W2STUTR <sup>2</sup> )	25,167	15,188	60.4	59.8
Second follow-up	F2 Student questionnaire (W4STUDENT <sup>2</sup> )	25,123	17,335	67.9	69.0
Second follow-up and base year	BY/F2 Student questionnaires (W4W1STU <sup>2</sup> )	25,123	15,909	62.5	63.3
	BY/F2 Student and BY Parent questionnaires <sup>9</sup> (W4W1STUP1 <sup>2</sup> )	25,123	12,888	50.1	51.3
Second follow-up, base year, and first follow-up	BY/F2 Student and BY/F1 Parent questionnaires <sup>5,10</sup> (W4W1STUP1P2 <sup>2</sup> )	11,927	5,427	44.6	45.5
Second follow-up, base year, first follow-up, and 2013 Update	BY/F1/U13/F2 Student questionnaires (W4W1W2W3STU <sup>2</sup> )	25,123	13,283	52.0	52.9
Postsecondary transcript	Postsecondary transcript (W5PSTRANS <sup>2</sup> )	17,338	13,160	71.2	75.9

See notes at end of table.

**Table 10. Number and percentage of completed surveys, high school transcript responses, postsecondary transcript and student records responses, or their combinations for the student sample, and associated recommended weights: PETS-SR—Continued**

Study round and high school transcript combinations	Data source(s) and recommended weights	Eligible	Participated	Weighted percent <sup>1</sup>	Unweighted percent
Postsecondary transcript, base year, first follow-up, and 2013 Update, second follow-up	Postsecondary transcript and BY/F1/U13/F2 Student questionnaires (W5W1W2W3W4PSTRANS <sup>2</sup> )	17,338	8,932	47.8	51.5
Postsecondary student records	Postsecondary student records (W5PSRECORDS <sup>2</sup> )	17,230	8,688	48.7	50.4
Postsecondary student records, base year, first follow-up, and 2013 Update, second follow-up	Postsecondary student records and BY/F1/U13/F2 Student questionnaires (W5W1W2W3W4PSRECORDS <sup>2</sup> )	17,230	5,936	32.8	34.4

<sup>1</sup> All weighted percentages are calculated using the student base weight, or the student base weight adjusted for subsampling or unknown eligibility, if applicable.

<sup>2</sup> Recommended weight, constructed to account for response to the data source.

<sup>3</sup> Recommended weight, not constructed specifically for response to the data source.

<sup>4</sup> Results for the math teacher questionnaire reflect students who were enrolled in a mathematics course in the base year; results for the science teacher questionnaire reflect students who were enrolled in a science course in the base year.

<sup>5</sup> Details of the parent subsample design are provided in section 3.3.4 of the *HSLs:09 Base Year to First Follow-up Data File Documentation* (Ingels et al. 2013).

<sup>6</sup> Participants are identified as sampled students who participated in both the base year and first follow-up and who have parent responses in both the base year and first follow-up.

<sup>7</sup> Only sampled students who participated in both the base year and 2013 Update with a responding base-year math teacher are considered participants.

<sup>8</sup> Only sampled students who participated in both the base year and 2013 Update with a responding base-year science teacher are considered participants.

<sup>9</sup> Only sampled students who participated in both the base year and second follow-up with a responding parent in the base year are considered participants.

<sup>10</sup> Only sampled students who participated in both the base year and second follow-up with a responding parent in the base year and first follow-up are considered participants.

NOTE: BY = base year; F1 = first follow-up; F2 = second follow-up; U13 = 2013 Update. All counts and computed rates are at the student level.

SOURCE: U.S. Department of Education, National Center for Education Statistics, High School Longitudinal Study of 2009 (HSLs:09) Postsecondary Education Transcript Study and Student Financial Aid Records Collection, Public-use Data File.

Choosing a weight for analyses can be complicated. To help in choosing a weight, researchers should first think in terms of the particular time period or data source of interest for the HSLs:09 population of students—base year, first follow-up, 2013 Update, High School Transcript, second follow-up, postsecondary transcript or student records, or some combination thereof. Next, researchers should consider the magnitude of nonresponse in the records included in the analyses and the associated nonresponse adjustment(s) for each weight.

As an example of how nonresponse magnitude might influence an analyst's decisions regarding which weight to use, consider a regression-based analysis. Records are excluded from a regression model if model covariates are missing, if the analysis

weight is zero, or both. Consider an example in which both parent *and* science teacher data are desired for a regression model to produce base-year student-level estimates. Using the rules above, two weights may be appropriate, W1PARENT and W1SCITCH. Both weights account for nonresponse in the respective contextual data sources (i.e., parent and science teacher nonresponse, respectively). However, because neither addresses nonresponse from both parents *and* science teachers, the use of either weight will be less than optimal. One approach is to conduct the regression analysis using both weights separately and if the conclusions of the analysis do not depend on the choice of weight, then report the conclusions using one of the weights. If the choice of weight produces different results, then another option is to select the weight which accounts for what is considered the most important source of nonresponse bias in the context of the analysis.

In the event that no weight accommodates interview data from all time periods and data sources of interest, researchers will have to assess the available weights to determine which weight should be used. A general rule of thumb is to select the weight that accounts for as many components of nonresponse as possible and, in the event of a tie, to select the weight that yields the most records for the analysis of interest. For illustration, suppose an analysis will use postsecondary transcript data, interview data from the 2013 Update, student interview data from the base year, and parent data from the base year. There is no analysis weight that explicitly accounts for parent nonresponse in the base year, student interview nonresponse in the base year, interview nonresponse in the 2013 Update, and missing postsecondary transcript data. However, there is one weight that accounts for three of the four sources of nonresponse, W5W1W2W3W4PSTRANS—this weight is therefore recommended for the analysis. Furthermore, whenever postsecondary transcript or student records data is used for analysis, then one of the four PETS-SR weights should be used. As discussed in section 5.3 and subsection 5.3.1.2, the PETS and SR weights represent subsets of the HSL:09 base-year population of students. The only two weights designed to represent the PETS subpopulation of students are W5PSTRANS and W5W1W2W3W4PSTRANS, and the only two weights designed to represent the SR subpopulation of students are W5PSRECORDS and W5W1W2W3W4PSRECORDS. Note that there is no specific weight constructed for students who are eligible for, and have, postsecondary transcripts and student records. If an analyst seeks to utilize postsecondary transcript and student records data together, then it is recommended to use one of the postsecondary student records weights, either W5PSRECORDS or W5W1W2W3W4PSRECORDS, as far fewer students have student records data than have transcript data.

*A note on incorporation of base-year teacher interview data into analyses.* Several additional elements of the study design speak to a need for caution in using the teacher data for

longitudinal analysis: (1) mathematics achievement was measured at the beginning of 9th grade and the end of 11th grade, but teacher characteristics were only measured for the fall of 9th grade; (2) teachers were not asked to rate or comment upon the individual HSLS:09 student; (3) very little curricular or classroom-level information was collected; and (4) students were linked to courses as represented by course titles (e.g., Algebra II, or Geometry) but not to a specific classroom that met at a specific time and place (e.g., Algebra II, section 3, meeting at 9 a.m.). These caveats should be kept in mind when dealing with the base-year teacher data.

As was stressed in the *HSLS:09 Base-Year Data File Documentation* (Ingels et al. 2011), the teacher sample does *not* constitute a nationally representative or school-representative sample of 9th-grade mathematics and science teachers. The two separate mathematics and science teacher samples were not independently selected but rather depend on a linkage to a sampled student who was selected for the study using probability methods and who both was enrolled in the requisite subject area and participated in the base year. Although it is possible to create teacher-level and course-level datasets using the base-year teacher data, they do not constitute valid generalizable probability samples of teachers. For this reason, neither a teacher ID nor statistical weights have been provided to support a teacher-level analysis. The teacher weights in the base year support use of teacher data only as an extension of the student record, with the student as the unit of analysis.

If base-year teacher data are used in conjunction with data from other time periods or from noninterview sources, the premise in selecting a weight as discussed above applies. Consider an example in which both first follow-up student data *and* base-year math teacher data are desired for a regression model to produce first follow-up student-level estimates. The likely weight for this analysis is W2STUDENT. This weight adjusts for the nonresponse associated with first follow-up student data but not for the nonresponse associated with base-year math teacher data. Researchers are encouraged to examine the pattern of missing data associated with the base-year teacher component and the W2STUDENT weight. If such an analysis suggests that the data are not necessarily missing at random, then experienced researchers may choose to investigate additional adjustments to the weights or to the data, such as an appropriate imputation model. Note, however, that the public-use file has limited information for use in such adjustments. Consequently, any subsequent adjustment could introduce more bias, not less, compared to using the data and weights in their published state.

## 5.5 Measures of Precision: Standard Errors and Design Effects

This section discusses the standard errors and design effects associated with HSLS:09 estimation. Readers may refer to appendix F for tables providing survey estimates, standard errors, and design effects for various domains of interest, computed using the primary postsecondary transcript weight.

### 5.5.1 Standard Errors

Complex sample designs, like that used for HSLS:09, result in data that violate the assumptions that are normally required to assess the statistical significance of results. The standard errors of the estimates from complex surveys may vary from those that would be expected if the sample were a simple random sample and the observations were independent. Some standard software packages, however, do not calculate standard error estimates that account for complex sampling designs used to select the school and student samples. This incorrect design assumption can lead to estimated variances and confidence intervals that are too small, which may lead to incorrect results from hypothesis tests. Variables have been created for HSLS:09 to support two methods of standard error estimation that account for the HSLS:09 complex sample design: (1) a BRR variance estimation method using the BRR weights and the associated analysis weight and (2) a linearization variance estimation method through a Taylor series approximation using analysis weights and variables that represent school sampling strata and primary sampling units.<sup>15</sup> Please note that variables to support these two methods of variance estimation are available to users of the restricted-use data, but only the BRR variance estimation method is supported for users of public-use data. Researchers are advised to use specialized software such as SUDAAN, SAS, or Stata that adjusts standard errors to account for the complex sampling design using one of these methods. Examples of code for these software programs are provided below.

The importance of correct variance estimation is further emphasized in this section through a discussion of the BRR and linearization methodologies.

The two methods of variance estimation supported through available HSLS:09 variables are BRR and Taylor series linearization. BRR variance estimation is supported with either the HSLS:09 restricted-use or public-use files. This method

---

<sup>15</sup> NCES statistical standards recommend the use of replicate variance estimation over linearization methods. The sample design variables, strata, and primary sampling units were suppressed from the public-use file as one measure of disclosure avoidance (see section 5.7 for information regarding the disclosure risk analysis and protection).

does not need the analysis stratum and primary sampling unit (PSU) identifiers but does require a large set of replicate weights along with the associated analysis weight. The replicate weights account for unequal selection probabilities, stratification, and clustering; incorporate unknown eligibility, nonresponse, and calibration adjustments; and produce standard error estimates that are in general slightly larger than the corresponding estimates calculated with linearization (Wolter 2007).

To create the BRR weights, the original school sampling strata were collapsed into 199 BRR strata with representation across the characteristics used in school sampling (i.e., school type, region, and locale) and two BRR PSUs were formed. The BRR strata were randomly assigned to rows of a  $200 \times 200$  Hadamard matrix containing a sequence of +1 and -1 values. The matrix is then used to assign certain cases a weight of 0 in order to form BRR base weights. The base weights were then adjusted using procedures similar to those implemented for the analysis weights.

The general formula for calculating a BRR variance estimate, used in software packages designed for survey estimation, is as follows:

$$var(\hat{\theta}) = \frac{1}{200} \sum_{a=1}^{200} (\hat{\theta}_{(a)} - \hat{\theta})^2 \quad (5-1)$$

where 200 is the number of HSLS:09 BRR weights,  $\hat{\theta}$  is the estimated value for a statistic of interest (e.g., mean) calculated with a particular analysis weight, and  $\hat{\theta}_{(a)}$  is the corresponding value calculated with the  $a$ th BRR (replicate) weight ( $a = 1, \dots, 200$ ).

Taylor series linearization variance estimation requires software that uses the analysis weight, analysis stratum, and PSU identifiers to compute standard errors that are adjusted to account for the complex sample design (see, e.g., Binder [1983]; Woodruff [1971]). The PSU and stratum identifiers are provided in two restricted-use variables, PSU and STRAT\_ID. The PSU variable contains a unique value randomly generated for each sampled school. The 450 values of STRAT\_ID were constructed in the base year by combining two to three schools into one analysis stratum in such a way as to maximize retention of the original two-stage sample design and also increase the precision of the estimates through the degrees of freedom (Chromy 1981). To lower disclosure risk, variables to support linearization variance estimation are only provided through the HSLS:09 restricted-use file, which, unlike the public-use file, contains the stratum and PSU variables.

Currently available software that can compute standard errors adjusted to account for a complex sample design includes SUDAAN,<sup>16</sup> SAS SURVEY procedures,<sup>17</sup> WesVar,<sup>18</sup> Stata,<sup>19</sup> R,<sup>20</sup> and SPSS.<sup>21</sup> Example SAS-callable SUDAAN code for producing estimated means and standard errors using the linearization and BRR methods are shown in figures 8 and 9, respectively. The corresponding Stata code is provided in figures 10 and 11, SAS code provided in figures 12 and 13, and R survey package code provided in figures 14 and 15. IBM SPSS code for the linearization method is provided in figure 16.

**Figure 8. Example SAS-callable SUDAAN code to calculate an estimated mean and linearization standard error for a postsecondary transcript student-level analysis**

```
PROC SORT DATA=<filename>;          *File sorted by nest variables;
  BY STRAT_ID PSU;
RUN;

PROC DESCRIPT DATA=<filename> DESIGN=WR;
  NEST STRAT_ID PSU / MISSUNIT;          *Analysis stratum/PSU;
  SUBPOPN (<domain variable = level>);   *Subset to reporting domain;
  WEIGHT W5PSTRANS;                      *Main analysis weight;
  VAR <analysis variable>;               *Analysis variable;
  PRINT MEAN SEMEAN / STYLE=NCHS;        *Mean and standard error;
RUN;
```

**Figure 9. Example SUDAAN code to calculate an estimated mean and replicate (BRR) standard error for a postsecondary transcript student-level analysis**

```
PROC DESCRIPT DATA=<filename> DESIGN=BRR;
  WEIGHT W5PSTRANS;                      *Main analysis weight;
  REPWGT W5PSTRANS001-W5PSTRANS200;      *BRR replicate weights;
  SUBPOPN (<domain variable = level>);   *Subset to reporting domain;
  VAR <analysis variable>;               *Analysis variable;
  PRINT MEAN SEMEAN / STYLE=NCHS;        *Mean and standard error;
RUN;
```

NOTE: BRR = balanced repeated replication.

<sup>16</sup> See <https://sudaansupport.rti.org/>.

<sup>17</sup> See the most recent *SAS User's Guide*, located at <https://support.sas.com/documentation/>.

<sup>18</sup> See <https://www.westat.com/capability/information-systems-software/wesvar>.

<sup>19</sup> See <https://www.stata.com/>.

<sup>20</sup> See <https://www.r-project.org/>.

<sup>21</sup> See <https://www.ibm.com/analytics/data-science/predictive-analytics/spss-statistical-software>.

**Figure 10. Example Stata code to calculate an estimated mean and linearization standard error for a postsecondary transcript student-level analysis**

```
SVYSET PSU [PWEIGHT=W5PSTRANS], STRATA (STRAT_ID) VCE(LINEAR),
singleunit(centered)

SVY, SUBP (<domain variable >) : MEAN < analysis variable >
```

**Figure 11. Example Stata code to calculate an estimated mean and replicate (BRR) standard error for a postsecondary transcript student-level analysis**

```
SVYSET [PWEIGHT=W5PSTRANS], BRRWEIGHT(W5PSTRANS001-
W5PSTRANS200) VCE(BRR) MSE

SVY, SUBP (<domain variable >) : MEAN < analysis variable >
```

NOTE: BRR = balanced repeated replication.

**Figure 12. Example SAS code to calculate an estimated mean and linearization standard error for a postsecondary transcript student-level analysis**

```
PROC SURVEYMEANS DATA=<filename> VARMETHOD=TAYLOR NOMCAR;
  STRATA STRAT_ID;                                *Analysis stratum;
  CLUSTER PSU;                                    *Analysis PSU;
  DOMAIN (<domain variable >);                    *Subset to reporting domain;
  WEIGHT W5PSTRANS;                                *Main analysis weight;
  VAR <analysis variable>;                          *Analysis variable;
RUN;
```

**Figure 13. Example SAS code to calculate an estimated mean and replicate (BRR) standard error for a postsecondary transcript student-level analysis**

```
PROC SURVEYMEANS DATA=<filename> VARMETHOD=BRR;
  WEIGHT W5PSTRANS;                                *Main analysis weight;
  REPWEIGHTS W5PSTRANS001-W5PSTRANS200;          *BRR replicate
  weights;
  DOMAIN (<domain variable >);                    *Subset to reporting domain;
  VAR <analysis variable>;                          *Analysis variable;
RUN;
```

NOTE: BRR = balanced repeated replication.

**Figure 14. Example R survey package code to calculate an estimated mean and linearization standard error for a postsecondary transcript student-level analysis**

```
mydesign<-svydesign(id=~PSU, strata=~STRAT_ID,
weights=~W5PSTRANS, data=mydata)
```

NOTE: For the R survey package (Lumley 2014), “mydesign” can be renamed to any name for an R object to hold the specification of the survey design, and “mydata” is the name of the current dataset.

**Figure 15. Example R survey package code to calculate an estimated mean and replicate (BRR) standard error for a postsecondary transcript student-level analysis**

```
mydesign<-svydesign(type="BRR", weights=~W5PSTRANS,
repweights="W5PSTRANS[001-200]",
combined.weights=FALSE, data=mydata)
```

NOTE: BRR = balanced repeated replication. For the R survey package (Lumley 2014), “mydesign” can be renamed to any name for an R object to hold the specification of the survey design, and “mydata” is the name of the current dataset.

**Figure 16. Example IBM SPSS complex samples code to calculate an estimated mean and linearization standard error for a postsecondary transcript student-level analysis**

```
CSPLAN ANALYSIS
/PLAN FILE='myfile.csaplan'
/PLANVARS ANALYSISWEIGHT=W5PSTRANS
/DESIGN STRATA=STRAT_ID CLUSTER=PSU
/ESTIMATOR TYPE=WR
```

NOTE: The name “myfile” should be replaced with the desired file name.

Standard errors for a select number of variables are provided in appendix F along with their design effects, which are discussed in the next section.

### 5.5.2 Design Effects

Design effects (*deff*) measure the relative efficiency of a sample design using particular items collected in the survey. These values are calculated as the ratio of two estimated variances,

$$deff = \frac{\hat{V}_d(\hat{\theta})}{\hat{V}_s(\hat{\theta})}, \quad (5-2)$$

for an estimated characteristic  $\hat{\theta}$ . The numerator value,  $\hat{V}_d(\hat{\theta})$ , is the estimated variance that properly accounts for the complex sample design and the variability associated with the analysis weights. The denominator value,  $\hat{V}_s(\hat{\theta})$ , is the estimated variance from a simple random sample (*srs*) design of the same sample size.

In addition to  $deff$ , the root design effect or  $deft$  may also be calculated. Like  $deff$ , this statistic also provides a measure of relative efficiency of a sample design but in terms of the standard errors:

$$deft = \sqrt{\frac{\hat{V}_d(\hat{\theta})}{\hat{V}_s(\hat{\theta})}}, \quad (5-3)$$

where the components are the same as defined for expression (5-2).

As noted in subsection 5.5.1, correct estimation of the variance of estimates requires the use of specialized software that can account for unequal selection probabilities, stratification, and clustering. In situations where software is unable to adjust for stratification and clustering but can accommodate weights, design effects may be used to approximate design-based variance and standard error estimates and thereby to produce associated test statistics that account for the estimated design-based variance.

The first step in approximating design-based variance estimates requires construction of normalized analysis weights. Given one of the analysis weights,  $w_i$ , defined in section 5.3, normalized analysis weights are defined as

$$w_{i,norm} = w_i * \frac{n}{\sum_{i=1}^n w_i} \quad (5-4)$$

where  $n$  corresponds to the number of observations with a positive weight,  $i$  indexes the set of respondents with a positive weight, and  $\sum_{i=1}^n w_i$  is the sum of the analysis weights.

There are three methods that may be used to produce  $t$  and  $F$  test statistics using approximated design-based variance estimates. The first method involves approximating the design-based variance estimate and using it to manually calculate the test statistics. In this first method, the normalized weights are used to estimate the simple random sampling variance or standard error of the estimator of interest using the available software. The design-based variance estimate may be approximated by multiplying the variance estimate produced from the software by an appropriate value of  $deff$ . Symbolically,

$$\hat{V}_d(\hat{\theta}) \sim \hat{V}_s(\hat{\theta}) * deff \quad (5-5)$$

where  $\hat{V}_s(\hat{\theta})$  is provided by the software, and  $deff$  may correspond to a specific estimate or may be the median<sup>22</sup> or mean of  $deff$  over several estimates. If the estimate of interest is for a subpopulation, then the value used for  $deff$  may be generated from a subgroup of respondents. The design effects reported in table 11 and those provided in appendix F may also be used for this second step. The approximate design-based variance estimates may be used to manually compute  $t$  and  $F$  test statistics.

The second method involves using the available software along with the normalized weights to generate  $t$  and  $F$  test statistics and then dividing the  $t$  statistic by an appropriate  $deff$  value and dividing the  $F$  statistics by an appropriate  $deff$  value.

The third method requires computing a new analysis weight by dividing the normalized weights by an appropriate value of  $deff$  and using this new analysis weight with the available software, using the test statistics produced with the software for inference.

- The HSLS:09 PETS-SR  $deff/deft$  analysis included two sets of variables: 16 variables from the postsecondary transcript data to assess design effects associated with the postsecondary transcript single-round weight, and a separate set of 19 variables from postsecondary student financial aid records data to assess design effects associated with the postsecondary student records single-round weight. As with the estimated standard errors, the  $deff$  and  $deft$  estimates were produced using final analysis weights and data that were edited, imputed (if applicable), and treated to limit disclosure risk. The  $deff$  estimates were calculated using a model-based formulation, corresponding to the  $deff4$  option in SUDAAN. As in the first follow-up, 2013 Update, and second follow-up, the item selection was informed by finding variables common to the HSLS:09 prior rounds' design effect analysis. Additionally, selected items utilized variables equivalent to those included in several other NCES studies involving postsecondary student and transcript data such as the ELS:2002, and the NPSAS:16. The  $deff$  and  $deft$  estimates are provided in appendix F for the 16 postsecondary transcript and 19 postsecondary student records items chosen using the above-specified criteria. The average  $deff$  and  $deft$  across both sets of items is presented in table 11. Design effects for key variables within the associated component for the weight were evaluated to ensure that the design effects are of acceptable size, such as achieving a  $deft$  below five for as many student domains as possible.

---

<sup>22</sup> Median design effects are provided in appendix F.

**Table 11. Average design effects (*deff*) and root design effects (*deft*) for postsecondary transcript and student records variables**

Characteristic <sup>1</sup>	Postsecondary transcript respondents	Final postsecondary transcript weight <sup>2</sup>		Postsecondary student records respondents	Final postsecondary student records weight <sup>3</sup>	
		Average <i>deff</i> <sup>4</sup>	Average <i>deft</i> <sup>5</sup>		Average <i>deff</i> <sup>4</sup>	Average <i>deft</i> <sup>5</sup>
<b>Total</b>	<b>13,160</b>	<b>6.2</b>	<b>2.4</b>	<b>8,688</b>	<b>4.7</b>	<b>2.1</b>
School type						
Public	10,010	5.5	2.3	6,518	4.2	2.0
Private	3,150	6.7	2.5	2,170	5.8	2.3
Region						
Northeast	2,213	5.1	2.2	1,305	5.5	2.3
Midwest	3,487	5.6	2.3	2,364	3.1	1.7
South	5,283	5.7	2.3	3,665	4.2	2.0
West	2,177	6.3	2.5	1,354	4.7	2.1
Locale						
City	3,982	8.7	2.8	2,716	6.6	2.5
Suburban	4,883	3.8	1.9	3,050	3.4	1.8
Town	1,474	5.3	2.3	1,043	4.6	2.1
Rural	2,821	4.7	2.1	1,879	3.7	1.9
Student sex						
Male	6,058	4.9	2.2	4,019	3.9	2.0
Female	7,102	4.9	2.2	4,669	4.0	2.0
Student race/ethnicity <sup>6</sup>						
Hispanic	1,872	4.5	2.1	1,165	3.8	1.9
Asian	1,245	4.2	2.0	785	6.5	2.5
Black	1,269	3.5	1.9	850	3.8	1.9
Other	8,774	4.3	2.0	5,888	3.4	1.8
Socioeconomic status (SES) <sup>7</sup>						
Low SES	1,394	4.0	2.0	868	4.7	2.1
Middle SES	7,444	4.5	2.1	4,876	3.7	1.9
High SES	4,322	4.7	2.1	2,944	3.9	1.9

<sup>1</sup> The school characteristics (school type, region, and locale) presented here reflect the information obtained during the HSLS:09 base year and do not contain updated information presented on the cumulative data file. The demographic characteristics (sex, race/ethnicity, and SES) presented here reflect information obtained during the HSLS:09 base year and updated in the first follow-up.

<sup>2</sup> Design effects computed using the W5PSTRANS weight.

<sup>3</sup> Design effects computed using the W5PSRECORDS weight.

<sup>4</sup> The formula for the design effect (*deff*) is provided in expression (5-2).

<sup>5</sup> The formula for the root design effect (*deft*) is provided in expression (5-3).

<sup>6</sup> Race/ethnicity as defined in the student questionnaire. Race categories exclude persons of Hispanic ethnicity.

<sup>7</sup> SES categories were defined using the SES quintile variable from the first follow-up (X2SESQ5), where X2SESQ5 = 1 (1st quintile) represents low SES, X2SESQ5 = 5 (5th quintile) represents high SES, and the three middle quintiles were classified as middle SES.

SOURCE: U.S. Department of Education, National Center for Education Statistics, High School Longitudinal Study of 2009 (HSLS:09) Postsecondary Education Transcript Study and Student Financial Aid Records Collection, Public-use Data File.

## 5.6 Unit and Item Nonresponse Bias Analysis

Unit and item nonresponse bias analyses are presented in this section, with unit nonresponse discussed in section 5.6.1 and item nonresponse discussed in section 5.6.2.

### 5.6.1 Unit Nonresponse Bias Analysis

NCES Statistical Standard 4-4-1 states that “Any survey stage of data collection with a unit or item response rate less than 85 percent must be evaluated for the potential magnitude of nonresponse bias before the data or any analysis using the data may be released. Estimates of survey characteristics for nonrespondents and respondents are required to assess the potential nonresponse bias” (Seastrom 2014).

The nonresponse bias in an estimated mean based on respondents  $\bar{y}_R$ , is the difference between the expected value of this mean and the target parameter,  $\pi$ , the population mean. Analysts can estimate the target parameter for variables that are observed for both respondents (R) and nonrespondents (NR) as follows:

$\hat{\pi} = (1 - \eta)\bar{y}_R + \eta\bar{y}_{NR}$ . In the equation,  $\hat{\pi}$  is the estimated population mean,  $\eta$  is the weighted unit (or item) nonresponse rate,  $\bar{y}_R$  is the observed weighted mean for respondents, and  $\bar{y}_{NR}$  is the weighted mean for nonrespondents. For variables that are from the frame rather than from the sample, analysts can estimate  $\pi$  without sampling error. They can then estimate bias as the difference between the respondent mean and the full-sample mean:  $\hat{B}(\bar{y}_R) = \bar{y}_R - \hat{\pi}$ . Equivalently, bias can be estimated as the difference between the mean for respondents and the mean for nonrespondents, multiplied by the weighted nonresponse rate:

$\hat{B}(\bar{y}_R) = \eta(\bar{y}_R - \bar{y}_{NR})$ . Relative bias provides a measure of the magnitude of the bias relative to the sample mean and is estimated as:  $\widehat{RB}(\bar{y}_R) = \hat{B}(\bar{y}_R) / \hat{\pi}$ .

Unit nonresponse bias analyses were conducted for the sets of respondents corresponding to the four analysis weights constructed for PETS-SR. Fifteen categorical variables were used to assess unit nonresponse bias. Several of the 15 variables are derived from sampling frame data and are not available in either restricted-use or public-use files. The 15 items are listed below. The items include 12 variables used in the nonresponse adjustment and 3 (Charter school status, Religious affiliation, and School is a regular secondary) that were not. Variable names are provided for those variables available in a restricted-use file.

- School type (X1CONTROL)
- 9th-grade enrollment percent by race

- Charter school status (A1SCHTYPE)
- Total school enrollment
- 9th-grade enrollment
- Number of full-time teachers (A1FTTCHRS)
- Student-to-teacher ratio
- Census region (X1REGION)
- School urbanicity (X1LOCALE)
- Range of grades in school (X1GRADESPAN)
- Religious affiliation of school
- School is a regular secondary
- Augmented sample-state (X1STATE)
- Sex (X2SEX)
- Race (X2RACE)

These 15 variables in total comprise 67 categories. The explicit categorization and category labels for each of the 15 items are provided in appendix D. For each category, estimates of bias were calculated and statistical significance tests conducted for each set of respondents corresponding to each of the four analysis weights.

The results of the nonresponse bias analyses to assess the potential reduction in bias attributable to base weight adjustments for nonresponse are described in the following sections, beginning with a description of the statistical tests for unit nonresponse bias (section 5.6.1.1).

#### **5.6.1.1 Test of nonresponse bias**

The VARGEN procedure in SUDAAN was used to estimate bias and conduct  $t$  tests to determine whether bias was significantly different from zero at a .05 level of significance. No multiple comparison adjustment was used in assessing the statistical significance of the tests of bias. Bias estimates were computed for each set of respondents associated with each of the four analysis weights. For each set of respondents, biases were estimated before nonresponse and calibration weight adjustments were applied to the sampling base weight adjusted for unknown eligibility and then estimated after nonresponse weight adjustments were applied to the sampling base weight adjusted for unknown eligibility. The base weight adjusted for unknown eligibility refers to the weight constructed in the unknown eligibility adjustments described in section 5.3.1.2. As was also discussed in section 5.3.1.2, PETS and SR had separate eligibility definitions, and thus there were two distinct base weights adjusted for the two different types of unknown eligibility. The weight adjusted for unknown eligibility with respect to PETS was the base weight adjusted for unknown eligibility for the unit nonresponse bias analysis of the two weights

concerning PETS (W5PSTRANS and W5W1W2W3W4PSTRANS). The weight adjusted for unknown eligibility with respect to SR was the base weight adjusted for unknown eligibility for the unit nonresponse bias analysis of the two weights concerning SR (W5PSRECORDS and W5W1W2W3W4PSRECORDS).

Table 12 contains a summary of the analysis for the four PETS-SR analysis weights; see appendix D for the detailed analysis tables. The results of these nonresponse bias analyses suggest that there is not a substantial bias on the variables examined due to nonresponse after adjusting for that nonresponse. However, it is not possible to directly assess bias on the transcript and SR data since these data are not available for nonrespondents.

**Table 12. Summary statistics for unit nonresponse bias analyses before and after weight adjustments for nonresponse, by HSLS:09 PETS-SR analysis weights: 2018**

Analysis weight	Significant bias tests at .05 level <sup>1</sup>		Absolute relative bias <sup>2</sup>		
	Percent before weight adjustment	Percent after weight adjustment	Median before weight adjustment	Median after weight adjustment	Percent relative change <sup>3</sup>
[W5PSTRANS] Postsecondary transcript	40.3	0.0	3.1	0.0	-100.0
[W5W1W2W3W4PSTRANS] Base year to first follow-up, to 2013 Update to second follow-up with postsecondary transcript	37.3	0.0	4.0	0.0	-100.0
[W5PSRECORDS] Postsecondary student records	37.3	0.0	5.2	0.0	-100.0
[W5W1W2W3W4PSRECORDS] Base year to first follow-up, to 2013 Update to second follow-up with postsecondary student records	44.8	0.0	5.8	0.0	-100.0

<sup>1</sup> "Before" and "after" are in reference to the nonresponse weight adjustment. A total of 67 statistical tests were performed; the number 67 was used as the basis for the reported percentages.

<sup>2</sup> The absolute relative bias is the absolute value of the (percent) relative bias where the percent relative bias is calculated as 100 multiplied by the estimated bias divided by the estimate computed using respondents and nonrespondents.

<sup>3</sup> The percent relative change is the percentage decrease in median absolute relative bias after weight adjustment. The formula for this was  $100 * (\text{median bias value after adjustment} - \text{median bias value before adjustment}) / \text{median bias value before adjustment}$ .

SOURCE: U.S. Department of Education, National Center for Education Statistics, High School Longitudinal Study of 2009 (HSLS:09) Postsecondary Education Transcript Study and Student Financial Aid Records Collection, Public-use Data File.

### 5.6.1.2 Postsecondary transcript student-level (W5PSTRANS) unit nonresponse bias analysis

In keeping with the NCES statistical standards, nonresponse bias analyses were performed for postsecondary transcript responses using the student analysis weight W5PSTRANS because, as shown in table 6, the weighted student response rate for the postsecondary transcript collection was 71.2 percent. Students for whom a

postsecondary transcript with sufficient information was collected from an institution the student attended were considered respondents for the purposes of postsecondary transcript collection. See section 5.2 for further details.

Approximately 40.3 percent of the 67 statistical tests conducted for the student-level unit response data identified bias statistically significant at the .05 significance level (see table 12) prior to adjusting the weights for nonresponse. After adjustment, no tests were statistically significant at the .05 level of significance, and the median absolute relative bias was reduced by 100.0 percent. Results of the 67 statistical tests are presented in table D-1 in appendix D. Additional comparisons between estimates produced after nonresponse adjustment and estimates produced after poststratification are provided in table D-2 in appendix D.

#### **5.6.1.3 *Base year to first follow-up, to 2013 Update to second follow-up with postsecondary transcript student-level (W5W1W2W3W4PSTRANS) unit nonresponse bias analysis***

As shown in table 6, the weighted unit response rate for the postsecondary transcript collection was 71.2 percent. However, the weighted unit response rate for students with responses in the postsecondary transcript collection with response in the base year, first follow-up, 2013 Update, *and* second follow-up was 47.8 percent.

Approximately 37.3 percent of the 67 statistical tests for this group of respondents identified statistically significant bias at the .05 significance level (see table 12) prior to adjusting the weights for nonresponse. After adjustment, no tests were statistically significant at the .05 level of significance, and the median absolute relative bias was reduced by 100.0 percent. The detailed analyses are shown in table D-3 in appendix D. Additional comparisons between estimates produced after nonresponse adjustment and estimates produced after poststratification are provided in table D-4 in appendix D.

#### **5.6.1.4 *Postsecondary student records student-level (W5PSRECORDS) unit nonresponse bias analysis***

In keeping with the NCES statistical standards, nonresponse bias analyses were performed for postsecondary student records responses using the student analysis weight W5PSRECORDS because, as shown in table 6, the weighted student response rate for the postsecondary student records collection was 48.7 percent. Students were considered respondents for the purposes of postsecondary student records collection if sufficient information regarding the student's financial aid was provided by the institution deemed to be the student's primary postsecondary institution. See section 5.1 for further details.

Approximately 37.3 percent of the 67 statistical tests conducted for the student-level unit response data identified bias statistically significant at the .05 significance level (see table 12) prior to adjusting the weights for nonresponse. After adjustment, no tests were statistically significant at the .05 level of significance, and the median absolute relative bias was reduced by 100.0 percent. Results of the 67 statistical tests are presented in table D-5 in appendix D. Additional comparisons between estimates produced after nonresponse adjustment and estimates produced after poststratification are provided in table D-6 in appendix D.

#### **5.6.1.5 Base year to first follow-up, to 2013 Update to second follow-up with postsecondary student records student-level (W5W1W2W3W4PSRECORDS) unit nonresponse bias analysis**

As shown in table 6, the weighted unit response rate for the postsecondary student records collection was 48.7 percent. However, the weighted unit response rate for students with responses in the postsecondary student records collection with response in the base year, first follow-up, 2013 Update, *and* second follow-up was 32.8 percent. Approximately 44.8 percent of the 67 statistical tests for this group of respondents identified statistically significant bias at the .05 significance level (see table 12) prior to adjusting the weights for nonresponse. After adjustment, no tests were statistically significant at the .05 level of significance, and the median absolute relative bias was reduced by 100.0 percent. The detailed analyses are shown in table D-7 in appendix D. Additional comparisons between estimates produced after nonresponse adjustment and estimates produced after poststratification are provided in table D-8 in appendix D.

### **5.6.2 Item Nonresponse Bias Analysis**

NCES Statistical Standard 4-4-3A states: “For an item with a low total response rate, respondents and nonrespondents can be compared on sampling frame and/or questionnaire variables for which data on respondents and nonrespondents are available. Base weights must be used in such analysis. Comparison items should have very high response rates. A full range of available items should be used for these comparisons. This approach may be limited to the extent that items available for respondents and nonrespondents may not be related to the low response rate item being analyzed” (Seastrom 2014).

Moreover, NCES Statistical Standard 1-3-5 states: “Item response rates (RRI) are calculated as the ratio of the number of respondents for whom an in-scope response was obtained ( $I^x$  for item  $x$ ) to the number of respondents who are asked to answer that item. The number asked to answer an item is the number of unit-level

respondents ( $I$ ) minus the number of respondents with a valid skip for item  $x$  ( $V^x$ ). When an abbreviated questionnaire is used to convert refusals, the eliminated questions are treated as item nonresponse. In the case of constructed variables, the numerator includes cases that have available data for the full set of items required to construct the variable, and the denominator includes all respondents eligible to respond to all items in the constructed variable” (Seastrom 2014). The item response rate is calculated as  $RRI^x = I^x / (I - V^x)$ .

All study items with a weighted response rate (weighted using either SR student analysis weight or the PETS student analysis weight<sup>23</sup>) of less than 85 percent were classified as having high item nonresponse and were included in the item nonresponse bias analyses. These variables and their response rates are described below in section 5.6.2.1.

The procedures for estimating and testing bias are the same as those used for unit nonresponse bias and are described in section 5.6.1. For each study item with less than an 85 percent response rate, as described above, bias estimates are computed by comparing item respondents to all other sample members who were eligible, or assumed eligible, for the item but did not respond to the item. NCES standards require that unit nonrespondents, whose item eligibility is unknown, must be assumed eligible for the item and must be treated as item nonrespondents. Consequently, bias estimates are computed using the student base weights since these weights are available for unit nonrespondents. The item nonresponse bias analysis was conducted using a subset of the frame variables used for the unit nonresponse bias analysis. The following school and student characteristics were available for both respondents and nonrespondents from the sampling frame and were used to assess item nonresponse bias:

- School type (X2CONTROL)
- Region of the United States (X2REGION)
- Locale (X2LOCALE)
- Sex (X2SEX)
- Race/ethnicity (X2RACE)

Item response rates are discussed in section 5.6.2.1 and results of the item nonresponse bias analysis are summarized in section 5.6.2.2. Detailed results for each item subject to nonresponse bias analysis appear in appendix tables D-9 through D-43.

---

<sup>23</sup> While analysis weights are used to construct item response rates, student base weights are used to carry out item-level nonresponse bias analyses.

### 5.6.2.1 Variables with high item nonresponse

All SR-PETS restricted-use student-level variables were reviewed to identify variables with a response rate below 85 percent. A total of 25 SR items and 10 PETS items had a response rate below 85 percent and were included in the nonresponse bias analysis. These variables and their response rates are given in tables 13 and 14.

Item response rates are calculated using both students for whom eligibility is known and students for whom eligibility is not known. Items that have a high completion rate among students with known eligibility may have a relatively small weighted response rate because students with unknown eligibility are assumed to be eligible and treated as nonrespondents. Similarly, items for which a low percentage of the population are eligible, may have a relatively high number of unknown eligible cases to known eligible. For example, the lowest weighted item response rate among SR items, 4.1 percent, was found for the 2011–2012 NSLDS variables *Deferred federal loans* (X5DEFER12), *Ever defaulted* (X5EVRDEF12), and *Federal loan entered forbearance* (X5FORBEAR12). Only 0.1 percent of students are known to be eligible for these items, while 4.1 of students were missing the information needed to match to NSLDS, and thus their eligibility for these items is unknown. The lowest item response rate among PETS items was 11.7 percent for *Remedial English courses: ratio of number known taken to known passed* (X5REMENRAT).

**Table 13. Student records items with a weighted item response rate below 85 percent using SR student weight (W5PSRECORDS)**

Variable name	Description	Percent of records by type of response			Unweighted item response rate	Weighted item response rate <sup>1</sup>
		Valid	Not applicable	Item missing		
X5PRIVLOANCUM	Student Records: Cumulative private (alternative) loans through June 30, 2016	79.0	.	21.0	79.0	77.5
X5PFYSEOGAMT	Student Records: Federal Supplemental Education Opportunity Grants at primary first year institution	70.8	.	29.2	70.8	72.9
X5PFYT4GRTAMT	Student Records: Total federal Title IV grants at primary first year institution	70.8	.	29.2	70.8	72.9
X5PFYNEEDAID	Student Records: Total need-based grants at primary first year institution	65.3	.	34.7	65.3	67.9
X5PFYTFEDWRK	Student Records: Federal work-study at primary first year institution	65.1	.	34.9	65.1	67.4
X5EVRFEDAPP	NSLDS: Applied for federal financial aid as of June 30, 2016	64.7	.	35.3	64.7	67.0
X5FEDAPP14	NSLDS: Applied for federal aid 2013–14	64.7	.	35.3	64.7	67.0
X5FEDAPP15	NSLDS: Applied for federal aid 2014–15	64.7	.	35.3	64.7	67.0
X5FEDAPP16	NSLDS: Applied for federal aid 2015–16	64.7	.	35.3	64.7	67.0
X5PFYCAMPAMT	Student Records: Federal campus-based aid (Perkins, SEOG, FWS) at primary first year institution	64.7	.	35.3	64.7	67.0
X5PFYFEDNEED	Student Records: Federal need-based aid at primary first year institution	64.7	.	35.3	64.7	67.0
X5PFYFEDPACK	Student Records: Federal Title IV aid package by type of aid at primary first year institution	64.7	.	35.3	64.7	67.0
X5PFYTITIVAIDREC	Student Records: Received any federal Title IV aid at primary first year institution	64.7	.	35.3	64.7	67.0
X5PFYTITIVAMT	Student Records: Total federal Title IV aid at primary first year institution	64.7	.	35.3	64.7	67.0

See notes at end of table.

**Table 13. Student records items with a weighted item response rate below 85 percent using SR student weight (W5PSRECORDS)—Continued**

Variable name	Description	Percent of records by type of response			Unweighted item response rate	Weighted item response rate <sup>1</sup>
		Valid	Not applicable	Item missing		
X5PFYNETPRICEGRT	Student Records: Tuition and fees minus Title IV, state, and institution grants at primary first year institution	63.8	.	36.2	63.8	66.1
X5PFYANYAIDREC	Student Records: Received any financial aid at primary first year institution	62.4	.	37.6	62.4	64.3
X5PFYNETPRICEALL	Student Records: Tuition and fees minus Title IV, state, and institution aid at primary first year institution	62.4	.	37.6	62.4	64.3
X5PFYPELLPACK	Student Records: Aid package with Pell Grants at primary first year institution	62.4	.	37.6	62.4	64.3
X5PFYTOTAID2	Student Records: Total federal (Title IV), state, and institutional aid at primary first year institution	62.4	.	37.6	62.4	64.3
X5DEFER13	NSLDS: Deferred federal loans 2012–13	1.4	94.5	4.1	25.3	33.0
X5EVRDEF13	NSLDS: Ever defaulted on a loan 2012–13	1.4	94.5	4.1	25.3	33.0
X5FORBEAR13	NSLDS: Federal loan entered forbearance 2012–13	1.4	94.5	4.1	25.3	33.0
X5DEFER12	NSLDS: Deferred federal loans 2011–12	0.1	95.8	4.1	1.9	4.1
X5EVRDEF12	NSLDS: Ever defaulted on a loan 2011–12	0.1	95.8	4.1	1.9	4.1
X5FORBEAR12	NSLDS: Federal loan entered forbearance 2011–12	0.1	95.8	4.1	1.9	4.1

<sup>1</sup> Weighted response rates were calculated with the SR student analysis weight (W5PSRECORDS).

NOTE: NSLDS = National Student Loan Data System; SR = Student Financial Aid Records.

SOURCE: U.S. Department of Education, National Center for Education Statistics, High School Longitudinal Study of 2009 (HSLS:09) Postsecondary Education Transcript Study and Student Financial Aid Records Collection, Restricted-use Data File.

**Table 14. Student records items with a weighted item response rate below 85 percent using PETS student weight (W5PSTRANS)**

Variable name	Description	Percent of records by type of response			Unweighted item response rate	Weighted item response rate <sup>1</sup>
		Valid	Not applicable	Item missing		
X5GPACRT	Transcript: GPA at first known certificate institution	3.2	96.3	0.5	85.4	83.7
X5GPALAST	Transcript: GPA at last known institution attended	84.2	1.8	14.0	85.8	83.7
X5HIGH10MAJ	Transcript: 10-category major of highest known degree obtained as of June 2016	9.5	88.0	2.5	79.0	79.7
X5HIGH11MAJ	Transcript: 11-category major of highest known degree obtained as of June 2016	9.5	88.0	2.5	79.0	79.7
X5HIGH23MAJ	Transcript: 23-category major of highest known degree obtained as of June 2016	9.5	88.0	2.5	79.0	79.7
X5HIGHCIP	Transcript: 6-digit CIP code of highest known degree obtained as of June 2016	9.5	88.0	2.5	79.0	79.7
X5STOPGT4M	Transcript: Count of known stopouts longer than 4 months	48.5	.	51.5	48.5	53.3
X5REMPSTRAT	Transcript: Remedial courses: ratio of number known taken to known passed	35.6	.	64.4	35.6	41.3
X5REMMTRAT	Transcript: Remedial mathematics courses: ratio of number known taken to known passed	28.0	.	72.0	28.0	32.9
X5REMENRAT	Transcript: Remedial English courses: ratio of number known taken to known passed	8.8	.	91.2	8.8	11.7

<sup>1</sup> Weighted response rates were calculated with the PETS student analysis weight (W5PSTRANS).

NOTE: GPA = grade point average; PETS = Postsecondary Education Transcript Study.

SOURCE: U.S. Department of Education, National Center for Education Statistics, High School Longitudinal Study of 2009 (HSL:09) Postsecondary Education Transcript Study and Student Financial Aid Records Collection, Restricted-use Data File.

### 5.6.2.2 Item nonresponse bias analysis results

Nonresponse bias results for each item listed in tables 13 and 14 are included in appendix D. For each item, bias was estimated and tested for each level of the five frame variables used, for a total of 16 estimates per item, as described in section 5.6.1.1.

Tables 15 and 16 summarize the bias ratios across all bias estimates. Bias ratios larger than 2.0 suggest the effect of item nonresponse may not be negligible. Of the 400 bias tests conducted across the 25 SR items, 38.8 percent had a bias ratio greater than 2.0. Of the 160 bias tests conducted across the 10 PETS items, 55.7 percent had a bias ratio greater than 2.0.

Table 17 and 18 summarize the significance tests and relative biases for all bias estimates. Overall, 40.3 percent of the bias estimates for SR items were statistically different from zero. The average relative bias is -1.6 and the median relative bias is -0.1. The average absolute relative bias is 18.2 and the median absolute relative bias is 7.0. The relative bias estimates varied a great deal by frame variable characteristic. For PETS items, 55.6 percent of the bias estimates were significantly different from zero. The average relative bias is -1.7 and the median relative bias is -1.6. The average absolute relative bias is 15.1 and the median absolute relative bias is 13.8.

Analysts should exercise caution when analyzing items where the results of the item nonresponse bias analysis suggest the presence of nontrivial levels of bias.

**Table 15. Frequency distribution of the estimated bias ratios for student records items**

Study instrument	Range of bias ratio <sup>1</sup>	Frequency <sup>2</sup>	Percent <sup>3</sup>
<b>Student records<sup>4</sup></b>	<b>Total</b>		<b>100.0</b>
	0 ≤ bias ratio < 2.0	245	61.3
	2.0 ≤ bias ratio < 5.0	117	29.3
	5.0 ≤ bias ratio	38	9.5

<sup>1</sup> The bias ratio is calculated as the estimated item nonresponse bias divided by the estimated standard error of the bias.

<sup>2</sup> The number of bias ratio calculations falling in the specified range of values.

<sup>3</sup> Percentage of bias ratio calculations falling in the specified range of values.

<sup>4</sup> The set of respondents used for bias estimation correspond to those in the student records (SR) data collection. Such students have a nonzero value for the SR student weight W5PSRECORDS.

NOTE: Detail may not sum to totals because of rounding.

SOURCE: U.S. Department of Education, Institute of Education Sciences, National Center for Education Statistics. High School Longitudinal Study of 2009 (HSL:09) Postsecondary Education Transcript Study and Student Financial Aid Records Collection, Restricted-use Data File.

**Table 16. Frequency distribution of the estimated bias ratios for transcript items**

Study instrument	Range of bias ratio <sup>1</sup>	Frequency <sup>2</sup>	Percent <sup>3</sup>
Student transcript <sup>4</sup>	<b>Total</b>		<b>100.0</b>
	0 ≤ bias ratio < 2.0	71	44.4
	2.0 ≤ bias ratio < 5.0	83	51.9
	5.0 ≤ bias ratio	6	3.8

<sup>1</sup> The bias ratio is calculated as the estimated item nonresponse bias divided by the estimated standard error of the bias.

<sup>2</sup> The number of bias ratio calculations falling in the specified range of values.

<sup>3</sup> Percentage of bias ratio calculations falling in the specified range of values.

<sup>4</sup> The set of respondents used for bias estimation correspond to those in the student transcript data collection. Such students have a nonzero value for the Postsecondary Education Transcript Study student analysis weight W5PSTRANS.

NOTE: Detail may not sum to totals because of rounding.

SOURCE: U.S. Department of Education, Institute of Education Sciences, National Center for Education Statistics. High School Longitudinal Study of 2009 (HSLs:09) Postsecondary Education Transcript Study and Student Financial Aid Records Collection, Restricted-use Data File.

**Table 17. Summary statistics for student records item nonresponse bias analyses using WSPSRECORDS weight**

School characteristics	Number of <i>t</i> tests	Percent <sup>1</sup> of significant <i>t</i> tests	Relative bias <sup>2</sup>		Absolute relative bias <sup>3</sup>	
			Average	Median	Average	Median
<b>Total</b>	<b>400</b>	<b>40.3</b>	<b>-1.6</b>	<b>-0.1</b>	<b>18.2</b>	<b>7.0</b>
School type						
Public	25	16.0	1.1	0.0	1.4	0.1
Private	25	16.0	-14.8	-0.0	17.4	1.0
Region						
Northeast	25	12.0	-21.5	-9.3	22.1	9.3
Midwest	25	28.0	-11.0	4.6	18.9	7.1
South	25	84.0	29.7	11.9	29.7	11.9
West	25	84.0	-15.3	-18.6	23.1	19.1
Locale						
City	25	0.0	3.1	-5.0	10.0	6.0
Suburban	25	16.0	-4.4	-2.3	6.4	2.7
Town	25	12.0	-2.3	10.6	20.2	15.0
Rural	25	0.0	1.0	0.9	2.2	1.3
Race/ethnicity <sup>4</sup>						
Hispanic	25	76.0	-9.7	-15.9	20.2	16.1
Asian	25	24.0	-21.4	-2.5	23.5	4.4
Black	25	24.0	42.4	10.2	43.6	10.2
Other	25	52.0	-4.9	3.0	9.7	3.5
Student sex						
Male	25	100.0	-20.6	-7.0	20.6	7.0
Female	25	100.0	22.3	6.5	22.3	6.5

<sup>1</sup> Percentage of *t* tests with  $p < 0.05$ .<sup>2</sup> Relative bias is calculated as 100 times the estimated bias divided by the weighted full-sample mean, using the student design weight adjusted for unknown eligibility.<sup>3</sup> Absolute relative bias is the absolute value of the relative bias.<sup>4</sup> Race/ethnicity as defined in the student questionnaire. Race categories exclude persons of Hispanic ethnicity.

SOURCE: U.S. Department of Education, National Center for Education Statistics, High School Longitudinal Study of 2009 (HSL:09) Postsecondary Education Transcript Study and Student Financial Aid Records Collection, Restricted-use Data File.

**Table 18. Summary statistics for student-level item nonresponse bias analyses using W5PSTRANS weight**

School characteristics	Number of <i>t</i> tests	Percent <sup>1</sup> of significant <i>t</i> tests	Relative bias <sup>2</sup>		Absolute relative bias <sup>3</sup>	
			Average	Median	Average	Median
<b>Total</b>	<b>160</b>	<b>55.6</b>	<b>-1.7</b>	<b>-1.6</b>		<b>15.1</b>
School type						
Public	10	60.0	1.2	0.6	1.5	1.2
Private	10	60.0	-15.2	-7.3	18.3	12.5
Region						
Northeast	10	30.0	-9.7	-1.9	9.7	1.9
Midwest	10	50.0	-10.9	-13.1	11.7	13.1
South	10	60.0	12.3	12.4	12.3	12.4
West	10	20.0	0.5	-5.5	11.5	11.0
Locale						
City	10	50.0	-17.2	-18.7	18.2	18.7
Suburban	10	20.0	-2.7	2.3	6.5	4.5
Town	10	70.0	28.6	30.8	28.6	30.8
Rural	10	10.0	8.0	7.3	10.8	9.5
Race/ethnicity <sup>4</sup>						
Hispanic	10	80.0	-3.8	-15.3	17.5	18.9
Asian	10	30.0	-20.3	-24.4	24.2	24.4
Black	10	90.0	-1.7	-16.8	34.5	36.8
Other	10	100.0	4.0	11.3	14.4	17.1
Student sex						
Male	10	80.0	-10.4	-12.3	10.8	12.3
Female	10	80.0	10.4	12.0	10.8	12.0

<sup>1</sup> Percent of *t* tests with  $p < 0.05$ .<sup>2</sup> Relative bias is calculated as 100 times the estimated bias divided by the weighted full-sample mean, using the student design weight adjusted for unknown eligibility.<sup>3</sup> Absolute relative bias is the absolute value of the relative bias.<sup>4</sup> Race/ethnicity as defined in the student questionnaire. Race categories exclude persons of Hispanic ethnicity.

SOURCE: U.S. Department of Education, National Center for Education Statistics, High School Longitudinal Study of 2009 (HSL:09) Postsecondary Education Transcript Study and Student Financial Aid Records Collection, Restricted-use Data File.

## 5.7 Single-value Item Imputation

Missing data in an otherwise complete study instrument occurs when a study respondent does not answer a particular question either intentionally (e.g., declined to answer a sensitive question) or unintentionally (e.g., missed one item within a set of related questions). Most statistical software packages exclude records that do not contain complete information. This is of great concern for multivariate analyses where a combination of missing values could greatly reduce the utility of the data.

To alleviate the problem of missing data from a respondent record, statistical imputation methods were employed for PETS-SR similar to those used for the HSLS:09 base year, first follow-up, 2013 Update, and second follow-up. Advantages of using imputed values include the ability to use all study respondent records in an analysis, which affords greater statistical power. Additionally, if the imputation procedure is effective (i.e., the imputed value is equal to, or close to, the true value), then the analysis results are possibly less biased than those produced with the incomplete data file.

A set of key analytic variables was identified for item imputation for study participants who responded to PETS-SR. Values were assigned in place of missing responses through single-value imputation for 21 student records variables. Indicator variables (flags) are included on the analysis file to allow users to easily identify the imputed values. The quality-control and evaluative procedures related to imputation are summarized in section 5.7.2.

### **5.7.1 *Imputed Items***

Twenty-one key analysis variables were identified for single-value imputation (see table 19) from the PETS-SR data. Additional variables were considered for this list but were excluded because of either high item-level response rates or they were deemed to be of lesser analytic importance.

**Table 19. Student records variables included in single-value imputation, by number and weighted percentage of values missing: 2018**

<b>Student records variables</b>	<b>Number of values imputed</b>	<b>Weighted percent imputed</b>
Institution merit-only grants (excludes athletic scholarships) at primary first year institution (R5YINSMERITNOATH)	480	5.5
State need-based only grants at primary first year institution (R5YSTNDONLY)	485	5.6
Institutional need-based grants at primary first year institution (R5YINSTNEED)	489	5.6
State merit-only grants at primary first year institution (R5YSTMERIT)	517	6.0
Institutional categorical grants at primary first year institution (R5YINSTCATGRT)	580	6.7
State grants based both on need and merit at primary first year institution (R5YSTNDMRT)	598	6.9
State need-based only grants at primary first year institution (R5YSTNOND1)	607	7.0
Athletic scholarships at primary first year institution (R5YINATHAMT)	617	7.1
State loans at primary first year institution (R5YSTLNAMT)	633	7.3
State work-study at primary first year institution (R5YSTWKAMT)	633	7.3
Institution military/armed forces grants at primary first year institution (R5YINSMILAMT)	634	7.3
Institution Veterans' education benefits at primary first year institution (R5YINSTVETAMT)	634	7.3
State military/armed forces grants at primary first year institution (R5YSTMILAMT)	634	7.3
State Veterans' education benefits at primary first year institution (R5YSTVETAMT)	634	7.3
Vocational rehabilitation and training at primary first year institution (R5YVOCHELP)	709	8.2
Institutional waivers, excludes employer waivers at primary first year institution (R5YINSWAIVNOEMP)	805	9.3
Institutional work-study at primary first year institution (R5YINSTWRK)	827	9.5
Institutional loans at primary first year institution (R5YINLNAMT)	829	9.5
Institutional tuition waivers for staff at primary first year institution (R5YEMPLWAIV)	835	9.6
Federal Supplemental Education Opportunity Grants at primary first year institution (R5YSEOGAMT)	2,533	29.2
Federal work-study at primary first year institution (R5YTFEDWRK)	3,036	35.0

SOURCE: U.S. Department of Education, National Center for Education Statistics, High School Longitudinal Study of 2009 (HSLs:09) Postsecondary Education Transcript Study and Student Financial Aid Records Collection, Public-use and Restricted-use Data Files.

Stochastic methods were used to impute the missing values for all variables included in table 19. Specifically, a weighted sequential hot-deck (WSHD) statistical imputation procedure (Cox 1980; Iannacchione 1982) was applied to the missing values for the variables in table 19 in the order in which they are listed. The WSHD procedure replaces missing data with valid data from a donor record (i.e., item respondent) within an imputation class. In general, variables with lower item nonresponse rates were imputed earlier in the process.

Imputation classes were identified using a recursive partitioning function in R. In addition to questionnaire items used to form the imputation classes, sorting variables

were used within each class to increase the chance of obtaining a close match between donor and recipient. If more than one sorting variable was chosen, a serpentine sort was performed where the direction of the sort—ascending or descending—changed each time the value of a variable changed. The serpentine sort minimized the change in the student characteristics every time one of the variables changed its value. With recursive partitioning, also known as a nonparametric classification tree or classification and regression tree (CART) analysis, the association of a set of questionnaire items and the variable requiring imputation is statistically tested (Breiman et al. 1984). The result is a set of imputation classes formed by the partition of the questionnaire items that are most predictive of the variable in question. The pattern of missing items within the imputation classes is expected to occur randomly so that the WSHD procedure can be used. The input questionnaire items included the sampling frame variables and variables imputed earlier in the ordered sequence or that were identified through skip patterns in the instrument and literature suggesting an association. The list of variables used as inputs to the CART procedure is provided in table G-1 of appendix G.

Cycling through the imputation variables, that is, the variables that will have imputed values, was part of the imputation process. Once the imputation variables are imputed the first time, the cycle returns and replaces the imputed values for the first imputation variable with the missing code. Then the imputation process re-imputes the first imputed variable using all variables, including the variables with imputed values, on the dataset. Next the imputation process moves to the second imputation variable, replaces the imputed values with missing values, and re-imputes the second variable. This process continues through all the imputation variables and is referred to as the second cycle. Five cycles were implemented for these imputation variables. The reasoning behind the use of cycling is that the imputed values will converge to a reasonable variable.

Finally, analysis weights were used to ensure that the population estimate calculated with data including the imputed values (post-imputation) did not change significantly from the estimate calculated prior to imputation (pre-imputation).

### **5.7.1.1 Imputation results**

Student records variables in table 19 are listed in the order in which they were imputed in addition to the number of values that were imputed for each variable. At each step, several quality-control procedures were used to maximize the utility of the imputed values. These are summarized in section 5.7.2.

### 5.7.2 *Evaluation of the Imputed Values*

After each value was imputed, a set of quality-control checks was implemented to ensure the highest quality of the imputed values. The unweighted distributions of the values before and after the imputation procedure were compared, both within and across the imputation classes, to identify large areas of change (see table G-2 of appendix G). Differences greater than 5 percent at the .05 significance level were flagged and examined to determine whether changes should be made to the imputation sort or class variables. Finally, data visualizations of value distributions before and after imputation were reviewed for potentially introduced bias.

The imputed variables' distributions within each imputation class were examined in order to identify classes where imputation might be done in a manner that does not emulate the raw data distribution. The visualization part is done for the variable in its entirety. Each variable is graphed three different ways—raw data, only imputed data, and raw plus imputed data—and compared for indications of introduced bias.

Multivariate consistency checks ensured that relationships among the imputation variables as well as between the imputation variables and key variables used for classification were maintained and that any special instructions for the imputation were implemented properly. For these checks, it was important to ensure that the imputation process did not create any new relationships that did not already exist in the observed data.

In any of the aforementioned checks, if there was any evidence of substantial deviation from the weighted sums or any identified inconsistencies, the imputation process was revised and rerun.

## 5.8 Disclosure Risk Analysis and Protections

Extensive confidentiality and data security procedures were employed for the PETS-SR data collection and data-processing activities. Data were prepared in accordance with NCES-approved disclosure-avoidance plans. The data disclosure guidelines were designed to minimize the likelihood of identifying individuals on the file by matching outliers or other unique data from external data sources. Because of the paramount importance of protecting the confidentiality of NCES data that contain information about specific individuals, data files were subject to various procedures to minimize disclosure risk. The PETS-SR data products and some of the disclosure treatment methods employed to produce them are described in the following sections. Details have been suppressed from this document to maintain the desired level of confidentiality.

### 5.8.1 *PETS-SR Data Products*

Data produced for the HSLS:09 PETS-SR data collection include restricted-use data and public-use data. Both the restricted- and public-use data include a student-level file. The student files contain responses and associated derived variables from the HSLS:09 PETS-SR data sources as well as all variables included in the base-year, first follow-up, 2013 Update and High School transcript, and second follow-up data files. Additional variables include those associated with survey-based analysis such as analysis strata and final analysis weights.

The disclosure treatment developed for the HSLS:09 PETS/SR data collection consisted of several steps:

- review of the collected data to identify items that may increase risk of disclosure;
- apply disclosure treatment to the high-risk items to lower the risk of disclosure;
- produce restricted-use data files that incorporate the disclosure-treated data; and
- produce public-use data files, constructed from the disclosure-treated restricted-use files, using additional disclosure limitation methods.

The disclosure treatment methods used to produce the PETS-SR data files include variable recoding, variable suppression, and swapping. These methods are described in section 5.8.2.

### 5.8.2 *Recoding, Suppression, and Swapping*

The disclosure treatment methods used to produce the PETS-SR data files include variable recoding, suppressing, and swapping. Some variables that had values with extremely low frequencies were recoded to ensure that the recoded values occurred with a reasonable frequency. Other variables were recoded from continuous to categorical values. In this way, rare events or characteristics have been masked for certain variables.

Some variables were classified as high risk and were suppressed from the public-use file. The suppressing techniques included removing the response from the file (i.e., reset to a “suppressed” reserve code) or removing records entirely from the public-use file.

Swapping was applied to certain items contained in the PETS-SR data files.

Swapping was implemented using NCES’ DataSwap software and utilized specific

and targeted, but undisclosed, swap rates. In data swapping, the values of the variables being swapped are exchanged between carefully selected pairs of records: a target record and a donor record. By doing so, even if an individual is tentatively identified, uncertainty remains about the accuracy and interpretation of the match because every record had some undisclosed probability of having been swapped.

Because perturbation (swapping) of the PETS-SR data could have changed the relationships between data items, an extensive data-quality check was carried out to assess and limit the impact of swapping on these relationships. For example, a set of utility measures for a variety of variables was evaluated pre- and post-treatment to verify that the swapping did not greatly affect the associations. Also, if the analysis determined that the components of a composite variable should be swapped, then the composite variable was reconstructed after swapping.

However, composite variables and their components could have been independently suppressed or recoded for inclusion in public-use files, resulting in a potential mismatch within the public-use file. In cases where recoding or suppression of composite variables and their components was carried out independently, public-use data users may not be able to recreate some of the composite variables provided in the public-use files. An example of this situation includes variables where the response categories have been collapsed for disclosure protection. The corresponding composite variable was derived from the full set of response categories as collected. Therefore, users who recalculate the composite variable with public-use information may see different results.

## Chapter 6. Data File Contents

This chapter provides an account of the file contents associated with HSLS:09 PETS-SR. Three types of products are available for researchers interested in using the HSLS:09 PETS-SR data: restricted-use files (NCES 2020-005) and public-use files (NCES 2020-021). Restricted-use files are restricted to users with a data-use license. These files include all source data at multiple levels—such as at the student level and at the student-by-institution level—and are typically acquired by users with relatively complex research questions. Data are available for download to all researchers through public-use files. These files allow for sophisticated or basic student-level analyses and are available across several common statistical packages. All products are updated versions of the HSLS:09 base-year through second follow-up data, meaning that the HSLS:09 PETS-SR data can be analyzed in conjunction with all previously released data for the HSLS:09 cohort.

### 6.1 PETS-SR Data Products

This section outlines how to access each data product and provides a detailed description of each product.

#### 6.1.1 *Restricted-use Data Products*

HSLS:09 PETS-SR restricted-use data are available on a DVD that includes restricted-use plain text data files and an electronic codebook (ECB) application. The data are available at no cost. A license is required to access the restricted-use data files. Details on obtaining a restricted-use license are available at <https://nces.ed.gov/statprog/instruct.asp>.

Because the restricted-use files provide data of all levels (e.g., student by institution), more advanced statistical analyses may require this data. The ECB application, an electronic version of a fully documented codebook, is easy to use and is designed to be accessible to researchers of all sophistication levels. It allows the user to browse all variables contained in the data files; search variable and value names for keywords of interest; review the question and item response wording; examine the definitions and logic used to develop composite variables; and export SAS, SPSS, or Stata syntax programs for statistical analyses. The ECB also displays the distribution and sample size for each variable. Analysts can use the ECB to export codebooks or generate

program code, including variable and value labels, in their desired programming language.

### 6.1.2 *Public-use Data Products*

The public-use data files include selected variables from the restricted-use files. Public-use data undergo more restrictive disclosure-avoidance treatment than the restricted-use data, including recoding and variable suppression as needed. The disclosure treatment developed for PETS-SR consisted of several steps:

- review of the collected data and identification of items that may increase risk of disclosure;
- application of disclosure treatment to the high-risk items to decrease the risk of disclosure;
- production of restricted-use data files that incorporate the disclosure-treated data; and
- production of public-use data files, constructed from the disclosure-treated restricted-use files, using additional disclosure limitation methods.

For more details on the disclosure treatment methods used to produce the HSLS:09 postsecondary transcripts and student records data files, please see section 5.8.

The public-use data are available via the web-based Online Codebook at <https://nces.ed.gov/onlinecodebook>. Online Codebook users can explore frequency distributions and select variables for download from the HSLS:09 public-use dataset. After a set of variables has been selected, the Online Codebook will also create a custom syntax file for use with the user's preferred software package (SAS, SPSS, Stata, R, or S-Plus). Alternatively, choosing a plain text file format (ASCII or CSV) allows for the data to be analyzed using most statistical programming languages.

## 6.2 Contents of the PETS-SR Data Products

The HSLS:09 PETS-SR restricted-use data contain the following data files:

- **Student data file (psstudent\_ruf).** Updated to include PETS-SR student-level composite variables
- **Weights file (psstudent\_brr\_ruf).** Updated to include PETS-SR student-level study and panel weights. See section 5.3 for more information about weight construction and section 5.4 for information on how to use the weights.
- Source data files
  - **Postsecondary institution data file (hsls\_pets\_institution).** Provides information about every postsecondary institution included in the postsecondary transcripts and student records data.
  - Postsecondary transcripts
    - **Student institution data (hsls\_pets\_stuinst).** Provides information for every institution the student attended (e.g., enrollment dates, awards received, total earned credits).
    - **Degree/major field of study data (hsls\_pets\_degmaj).** Provides information for every degree or major field of study indicated on the student transcripts (e.g., degree program, major CIP code, date received (if applicable), honors).
    - **Term data (hsls\_pets\_term).** Provides information for every term indicated on the student transcripts (e.g., term dates, honors, earned credits, GPA).
    - **Test data (hsls\_pets\_test).** Provides information for every test indicated on the student transcripts (e.g., test name, score, date).
    - **Course data (hsls\_pets\_course).** Provides information for every course indicated on the student transcripts (e.g., course name, CCM code, credits earned, and grade received).
  - Postsecondary SR files
    - **Student institution file (hsls\_sr\_stuinst).** Includes data for each student-institution pair, including demographics, standardized test scores, yearly enrollment flags, etc.
    - **Student institution by year file (hsls\_sr\_stuinstyr).** Includes data for each academic year a student was enrolled (or potentially

enrolled) at the given institution, including GPA, student budget, financial aid flags, etc.

- **Degree/major field of study file (hsls\_sr\_degmaaj).** Provides yearly information for every degree and major field of study, including CIP codes and required credit hours.
  - **Term file (hsls\_sr\_terms).** Lists every term in which the student was enrolled (or potentially enrolled) with enrollment status and enrolled credit hours.
  - **Test file (hsls\_sr\_test).** Lists any reported SAT and ACT scores.
  - **Financial aid award file (hsls\_sr\_aid).** Lists all financial aid awarded by source and program type.
- o NSLDS data. All files include data through January 2018.
- **Federal grant file (hsls\_nsllds\_pell).** Includes complete award histories, amounts, and pertinent dates for federal grants such as Pell Grants and National Science and Mathematics Access to Retain Talent (SMART) Grants.
  - **Federal loan file (hsls\_nsllds\_loan).** Lists information on federal loans borrowed, such as the loan program, status, and pertinent dates.
  - **Award origin file (hsls\_nsllds\_award).** Lists information on federal loans awarded, such as the year, dependency status, start date, and end date.
  - **Non-Stafford loan default file (hsls\_nsllds\_defnonstaf).** Provides statuses and start and end dates for default occurrences on non-Stafford Loans.
  - **Stafford loan default file (hsls\_nsllds\_defstaf).** Provides statuses and start and end dates for default occurrences on Stafford Loans.
  - **Enrollment Status file (hsls\_nsllds\_enroll).** Provides enrollment status codes, effective date, and credential level of program.
  - **Federal loan deferment file (hsls\_nsllds\_loandefer).** Includes information for each deferment-period update, including the type of deferment, start date, and end date.
  - **Federal loan delinquency file (hsls\_nsllds\_loandelinq).** Includes information for each delinquency period, including the beginning date and end date.

- **Federal loan disbursement file (hsls\_nsls\_loandis).** Includes information for each disbursement, including the disbursement date, and disbursement amount
- **Federal loan forbearance file (hsls\_nsls\_loanforbear).** Includes information for each forbearance period, including the forbearance type, beginning date, and end date.
- **Outstanding interest balance history file (hsls\_nsls\_oib).** Provides a history of each loan's outstanding interest balance.
- **Outstanding principal balance history file (hsls\_nsls\_opb).** Provides a history of each loan's outstanding principal balance.
- **Federal loan repayment file (hsls\_nsls\_rpmtplan).** Details each loan's repayment plan over time, including the type, monthly payment amount, and pertinent dates.

Table 20 provides an indication of which ECB files exist as a public-use file and which files are new as of the PETS-SR release.

**Table 20. PETS-SR data products: 2018**

ECB display order	File name	File description	PUF version exists	New this release
1	psstudent_ruf	Student File	X	Added X5 variables
2	school	School File	X	
3	sch_hstrns	HS Transcript School File		
4	stu_hstrns	HS Transcript Student School File		
5	sch_course	HS Transcript School Course File		
6	stu_course	HS Transcript Student Course File		
7	f2stu_inst	Student-Institution File		
8	f2stu_inst_prog	Student-Institution-Program File		
9	cps1314	CPS 2013-14 File		
10	cps1415	CPS 2014-15 File		
11	cps1516	CPS 2015-16 File		
12	cps1617	CPS 2016-17 File		
13	cps1718	CPS 2017-18 File		Replacement file
14	cps1819	CPS 2018-19 File		X
15	cps1920	CPS 2019-20 Preliminary File		X
16	hsls_pets_institution	PETS Institution File		X
17	hsls_pets_stuinst	PETS Student-Institution File		X
18	hsls_pets_degmaaj	PETS Degree Major file		X
19	hsls_pets_term	PETS Term File		X
20	hsls_pets_test	PETS Test File		X
21	hsls_pets_course	PETS Course File		X
22	hsls_sr_stuinst	SR Student-Institution File		X
23	hsls_sr_stuinstyr	SR Student-Institution by Year File		X
24	hsls_sr_degmaaj	SR Degree Major File		X
25	hsls_sr_terms	SR Terms File		X
26	hsls_sr_test	SR Test File		X
27	hsls_sr_aid	SR Aid File		X
28	hsls_nsls_pell	NSLDS Pell Grant File		Replacement file
29	hsls_nsls_loan	NSLDS Loan File		Replacement file
30	hsls_nsls_award	NSLDS Award Origin File		X
31	hsls_nsls_defnonstaf	NSLDS Non-Stafford Loan Default File		X
32	hsls_nsls_defstaf	NSLDS Stafford Loan Default File		X
33	hsls_nsls_enroll	NSLDS Enrollment Status File		X
34	hsls_nsls_loandefer	NSLDS Loan Deferment File		X
35	hsls_nsls_loandelinq	NSLDS Loan Delinquency File		X
36	hsls_nsls_loandis	NSLDS Loan Disbursement File		X
37	hsls_nsls_loanforbear	NSLDS Loan Forbearance File		X
38	hsls_nsls_oib	NSLDS Outstanding Interest Balance File		X
39	hsls_nsls_opb	NSLDS Outstanding Principal Balance File		X
40	hsls_nsls_rpmptplan	NSLDS Loan Repayment Plan File		X
41	psstudent_brr_ruf	Student BRR File		Added W5 variables

NOTE: BRR = balanced repeated replication; CPS = Central Processing System; ECB = electronic codebook, HS = high school; NSLDS = National Student Loan Data System; PETS = Postsecondary Education Transcript Study; PUF = public-use file; SR = Student Financial Aid Records.

SOURCE: U.S. Department of Education, National Center for Education Statistics, High School Longitudinal Study of 2009 (HSLS:09) Postsecondary Education Transcript Study and Student Financial Aid Records Collection.

## 6.3 Variable Naming Schema

All variable names include a prefix to help users easily identify the source of the data used in the variable, the round in which the data were collected, and the appropriate file or level at which the data are reported. Variable prefixes adhere to the following convention: the first character indicates the data source, the second character indicates the study round, the third character indicates the data level or file, and the remainder is a descriptive name that identifies the information captured by the variable.

The following first characters are associated with the PETS-SR data:

- X—Composite variables,
- W—Weights,
- T—Transcripts,
- R—Student records, and
- I—Institution.

The second character (study round indicator) is a “5” for all PETS-SR variables.

The following third characters are included on source data files for PETS-SR:

- A—Student aid,
- C—Courses,
- D—Degree/major field of study,
- M—Terms,
- S—Student institution,
- X—Tests, and
- Y—Student institution year.

As an example, the variable T5SHIGHAWD (Highest award at the institution) indicates that the data are from transcripts (T), collected during the HSLs:09 PETS-SR collection (5), and reported on the student institution file as a student-by-institution-level variable (S). Appendix H provides a listing of all PETS-SR variables, including the file name, variable name, and variable label for the subset of new data added to the HSLs:09 restricted-use files.

## 6.4 Missing Data

As mentioned in section 4.2, when data are missing, negative integers called *reserve codes* are inserted to indicate the cause of the missing data. For example, reserve codes allow distinctions to be drawn between an unknown value and a value that

does not apply to a sample member. These codes are used to specify missing data at the item and unit level across data files. Not only do these codes help delineate missing data within submitted postsecondary transcripts and student financial aid records, but records are added to hold a place for nonrespondents as well.

### 6.4.1 Reserve Codes

Table 21 provides a listing of the reserve code values employed across the PETS-SR data files.

**Table 21. Reserve code values: 2018**

Value	Description
-1	<b>Item missing, don't know</b> Used when a respondent indicated "Don't know" as a response.
-2	<b>Placeholder record</b> Used to hold a place for a record with an unknown amount of missing data. For example, if no academic terms are reported, one placeholder record is included in the term file, though the number of terms in which the student was enrolled is unknown.
-3	<b>Implied "No" or zero</b> Used when the item was left blank by the respondent, but based on other responses, the missing value is implied to be a "No" or a zero.
-4	<b>Item missing, unable to determine applicability</b> Used for a nested item when the associated gate item is left blank.
-5	<b>Data suppressed</b> Used on the student-level and school-level public-use data files to suppress data.
-6	<b>Out of range</b> Used when the value reported by the institution was outside the valid range for that field.
-7	<b>Item missing, not applicable</b> Used for questions that are not applicable based on information already known from a prior answer or another data source.
-8	<b>Unit missing</b> Used for all student-level variables when a sample member is a nonrespondent to either PETS or SR.
-9	<b>Item missing, response not provided</b> Used for questions that are not answered within a survey when the respondent was eligible for the question.

SOURCE: U.S. Department of Education, National Center for Education Statistics, High School Longitudinal Study of 2009 (HSL:09) Postsecondary Education Transcript Study and Student Financial Aid Records Collection.

### 6.4.2 Placeholder Records

Records were added to PETS-SR source files to hold a place for nonrespondents and unknown quantities of missing data. Specifically, on PETS files, if no data were reported for terms, courses, or degrees/majors, the associated file includes a single placeholder record to indicate that the data were not submitted by the institution, hence the number of records missing is unknown. Similarly, for the SR files, if no

data were reported for student institution academic years, terms, degrees/majors, or financial aid awards, one placeholder record is included on the respective file.

## 6.5 Composite Variables

A set of composite variables—also called *derived variables*—has been created for each round of HSLS:09, and a new set has been added to the student data file that incorporate PETS, SR, and NSLDS data. The new composite variables are generated with responses from two or more source variables, potentially from multiple data sources. Composite variable descriptions may be found in appendix I. The HSLS:09 second follow-up data products inherit composite variables from prior rounds as well as those newly created with data from the PETS-SR collection.

Most of the composite variables can be used as classification variables or independent variables in data analysis. Some of the key SR composites have undergone imputation to address missing responses. Note that all imputed versions of variables have been flagged. Variables with imputed data have a separate imputation flag variable with similar naming convention (\*\_IM suffix), and that imputation flag variable indicates which cases are imputed and the source of the imputation. For example, X5PFYTFEDWRK\_IM=2 where data are imputed for X5PFYTFEDWRK. For more on the imputation process, see section 5.7.

One of the goals of the student records composite variables is to report on the initial aid package students received when transitioning from high school to college. To do so, project staff used student financial aid records and postsecondary transcript data to identify a student's first academic year enrolled post-high school and the institution at which the student was primarily enrolled. All composite variables associated with the first primary institution record are denoted with a prefix of "X5PFY."

To the extent possible, the first primary postsecondary record aligns with the manner whereby the first "real" postsecondary institution was defined in ELS:2002, to maintain consistency across the secondary longitudinal studies program. Specifically, the first primary postsecondary record is generally the institution and associated academic year with the earliest start date post-high school where an academic year is defined to be July 1 through June 30. An exception was made for the first institution, if (1) enrollment at the first chronological institution was during the summer (i.e., the enrollment begins in May, June, or July and ends by August); (2) the summer enrollment was in the same calendar year and follows high school completion/exit; and, (3) fall enrollment (i.e., the enrollment began in August, September, or October)

was observed at a different postsecondary institution of that same calendar year, following high school completion/exit. If all the above conditions are met, the next institution with the next earliest start date was selected. This exception was made in order to exclude summer enrollments immediately following high school completion/exit and immediately preceding fall enrollment at another institution. If there was any substantive gap between high school and summer enrollment, or summer enrollment and the next enrollment spell, no exception was made. Additionally, if enrollment was observed in both summer and fall at the same institution following high school, no exception was made.

An exception was made for the first academic year if the (1) first enrolled month was May or June and (2) enrollment was observed in September or October of the same calendar year at the same institution. In this case, the academic year with the September/October enrollment was selected. Note that this exception was made in order to exclude an academic year if the only enrollment observed was in the last two months of the academic year and that enrollment was immediately followed by fall enrollment.

## 6.6 Data Anomalies and Considerations

The variables X2MTHINT and X4EVRTRANSHS had data errors in the prior release. Thus, the data have been corrected and the variables renamed as X2MTHINT\_R and X4EVRTRANSHS\_R. Details regarding the errors are as follows:

- X2MTHINT “Scale of student’s interest in fall 2009 math course” – One of the inputs to this composite variable had been erroneously programmed to indicate students whose favorite subjects were “science” (S2FAVSUBJ=3) instead of “math” (S2FAVSUBJ=6).
- X4EVRTRANSHS “Ever transferred from base year high school” – This composite was erroneously programmed to indicate students who were known to have transferred prior to the F2 round as transferred. Any student who was known to not transfer prior to the F2 round had to be updated as X2EVRTRANSHS\_R=0.

The following variables had inaccurate descriptions in the prior release. The variable descriptions have been corrected as:

- S4PRE\_03 – This variable was loaded into the second follow-up instrument for each sample member prior to survey administration. Second follow-up

respondents for whom transcript data indicate whether a student transferred high schools or not were assigned a value of “1” and were not administered S4TRANSFERHS. All other second follow-up respondents were assigned a value of “0” and were not administered S4TRANSFERHS.

- S4PRE\_04 – This variable was loaded into the second follow-up instrument for each sample member prior to survey administration. Second follow-up respondents for whom 2013 Update data indicated applying to college (i.e.,  $(S3CLGID < 0 \text{ and } S3CLGAPPID1 < 0 \text{ and } S3CLGAPPID2 < 0)$  or  $(S3CLGID > 0 \text{ and } S3CLGAPPID1 = -4 \text{ and } S3CLGAPPID2 = -4)$ ) were assigned a value of “1” and were administered college application questions in section B. All other second follow-up respondents were assigned a value of “0” and were not administered college application questions.
- X1SCHASIAN – Changed reference of A1ASIANSTU to A1ASIANPISTU in the description.
- X1 and X2 scales – Changed reference of “The coefficient of reliability (alpha) for the scale is .65” to “The coefficient of reliability (alpha) for the scale is .65 or higher” in the descriptions.

## References

- Adelman, C. (2006). *The Tool Box Revisited: Paths to Degree Completion from High School through College*. Washington, DC: U.S. Department of Education.
- Binder, D.A. (1983). On the variances of asymptotically normal estimators from complex surveys. *International Statistical Review*, 51(3), 279–292.
- Bozick, R., and Ingels, S.J. (2008). *Mathematics Coursetaking and Achievement at the End of High School* (NCES 2008-319). U.S. Department of Education. Washington, DC: National Center for Education Statistics, Institute of Education Sciences.
- Breiman, L., Friedmand, J., Olshen, R., and Stone, C. (1984). *Classification and Regression Trees*. Belmont, CA: Wadsworth, Inc.
- Bryan, M. and Simone, S. (2012). *2010 College Course Map* (NCES 2012-162REV). U.S. Department of Education. Washington, DC: National Center for Education Statistics, Institute of Education Sciences.
- Chaney, B., Burgdorf, K., and Atash, N. (1997). Influencing Achievement Through High School Graduation Requirements. *Educational Evaluation and Policy Analysis*, 19(3): 229–244.
- Chen, X. (2009). *Students Who Study Science, Technology, Engineering, and Mathematics (STEM) in Postsecondary Education* (NCES 2009-161). U.S. Department of Education. Washington, DC: National Center for Education Statistics, Institute of Education Sciences.
- Chromy, J.R. (1981). Variance Estimators for a Sequential Sample Selection Procedure. In D. Krewski, R. Platek, and J.N.K. Rao (Eds.), *Current Topics in Survey Sampling* (pp. 329–347). New York: Academic Press.
- Cohen, J. (1960). A Coefficient of Agreement for Nominal Scales. *Educational and Psychological Measurement*, 20(1): 37–46.
- Cool, V.A., and Keith, T.Z. (1991). Testing a Model of School Learning: Direct and Indirect Effects on Academic Achievement. *Contemporary Educational Psychology*, 16(1): 28–44.

- Cox, B.G. (1980). The Weighted Sequential Hot Deck Imputation Procedure. *Proceedings of the Section on Survey Research Methods* (pp. 721–726). Alexandria, VA: The American Statistical Association.
- Duprey, M.A., Pratt, D.J., Jewell, D.M., Cominole, M.B., Fritch, L.B., Ritchie, E.A., Rogers, J.E., Wescott, J.D., Wilson, D.H. (2018). *High School Longitudinal Study of 2009 (HSLs:09) Base-Year to Second Follow-Up Data File Documentation* (NCES 2018-140). U.S. Department of Education. Washington, DC: National Center for Education Statistics, Institute of Education Sciences.
- Folsom, R.E and Singh, A.C. (2000). The Generalized Exponential Model for Sampling Weight Calibration for Extreme Values, Nonresponse, and Poststratification. *Proceedings of the Section on Survey Research Methods of the American Statistical Association*, (pp. 598–603). Alexandria, VA: American Statistical Association.
- Iannacchione, V.G. (1982). Weighted Sequential Hot Deck Imputation Macros. In *Proceedings of the Seventh Annual SAS Users Group International Conference*, 759–763.
- Ingels, S.J., Pratt, D.J., Herget, D.R., Burns, L.J., Dever, J.A., Ottem, R., Rogers, J.E., Jin, Y., and Leinwand, S. (2011). *High School Longitudinal Study of 2009 (HSLs:09): Base-Year Data File Documentation* (NCES 2011-328). U.S. Department of Education. Washington, DC: National Center for Education Statistics, Institute of Education Sciences.
- Ingels, S.J., Pratt, D.J., Herget, D.R., Dever, J.A., Fritch, L.B., Ottem, R., Rogers, J.E., Kitmitto, S., and Leinwand, S. (2013). *High School Longitudinal Study of 2009 (HSLs:09) Base Year to First Follow-Up Data File Documentation* (NCES 2014-361). U.S. Department of Education. Washington, DC: National Center for Education Statistics, Institute of Education Sciences.
- Ingels, S.J., Pratt, D.J., Herget, D., Bryan, M., Fritch, L.B., Ottem, R., Rogers, J.E., and Wilson, D. (2015). *High School Longitudinal Study of 2009 (HSLs:09) 2013 Update and High School Transcript Data File Documentation* (NCES 2015-036). U.S. Department of Education. Washington, DC: National Center for Education Statistics, Institute of Education Sciences.
- Meyer, R.H. (1998). *The Production of Mathematics Skills in High School*. Chicago and Madison, WI: Harris School of Public Policy Studies, The University of Chicago, and Wisconsin Center for Education Research.

- Rock, D.A., and Pollack, J.M. (1995). *Mathematics Coursetaking and Gains in Mathematics Achievement* (NCES 95-714). U.S. Department of Education. Washington, DC: National Center for Education Statistics.
- Seastrom, M.M. (2014). *2012 Revision of NCES Statistical Standards*. (NCES 2014-097). U.S. Department of Education. Washington, DC: National Center for Education Statistics, Institute of Education Sciences.
- StandardsWork. (2006). *12th Grade NAEP Revised: What Stakeholders Say*. Prepared for National Assessment Governing Board, National Assessment of Educational Progress. Washington, DC: Author.
- Wolter, K. (2007). *Introduction to Variance Estimation*. Second Edition. New York: Springer-Verlag.
- Woodruff, R.S. (1971). A Simple Method for Approximating the Variance of a Complicated Estimate. *Journal of the American Statistical Association*, 66, 411–414.